

# HIPEAC

COMPILATION ARCHITECTURE

**HIGH PERFORMANCE AND EMBEDDED ARCHITECTURE AND COMPILATION**



## THE HIPEAC VISION FOR ADVANCED COMPUTING IN HORIZON 2020



**M. Duranton, D. Black-Schaffer,  
K. De Bosschere, J. Maebe**

This document was produced as a deliverable of the FP7 HiPEAC Network of Excellence under grant agreement 287759. March 2013

The editorial board is indebted to Dr Max Lemke and to Dr Panos Tsarchopoulos of the Complex Systems and Advanced Computing unit of the Directorate-General for Communications Networks, Content and Technology of the European Commission for their active support to this work.

# CONTENTS

<b>CONTENTS</b>	<b>1</b>	<b>3. TECHNOLOGY CONSTRAINTS AND OPPORTUNITIES</b>	<b>27</b>
<b>EXECUTIVE SUMMARY</b>	<b>3</b>	<b>3.1. CONSTRAINTS</b>	<b>27</b>
<b>PREFACE</b>	<b>6</b>	3.1.1. HIGH PERFORMANCE HARDWARE BLOCKED BY THE FOUNDRY COSTS	28
<b>1. ADVANCED COMPUTING SYSTEMS RECOMMENDATIONS FOR HORIZON 2020</b>	<b>7</b>	3.1.2. THE POWER CONSUMPTION COST	28
<b>1.1. STRATEGIC AREA 1: EMBEDDED SYSTEMS</b>	<b>8</b>	3.1.3. COMPLEXITY: THE ACHILLES HEEL OF SOFTWARE	31
1.1.1. COST-EFFECTIVE DESIGN OF EMBEDDED SYSTEMS	9	<b>3.2. OPPORTUNITIES</b>	<b>31</b>
1.1.2. COST-EFFECTIVE CERTIFICATION	9	3.2.1. ENTERING THE THIRD DIMENSION	31
1.1.3. SECURE EMBEDDED SYSTEMS	9	3.2.2. SILICON PHOTONICS	33
<b>1.2. STRATEGIC AREA 2: DATA CENTER COMPUTING</b>	<b>9</b>	3.2.3. WIRELESS CONNECTIVITY	33
1.2.1. LOW-POWER MICRO SERVERS AND MODULES	9	3.2.4. EMERGING MEMORY TECHNOLOGIES	33
1.2.2. NETWORK AND STORAGE I/O	10	3.2.5. STOCHASTIC/APPROXIMATE COMPUTING	35
<b>1.3. STRATEGIC AREA 3: MOBILE SYSTEMS</b>	<b>10</b>	3.2.6. NOVEL ARCHITECTURES	36
1.3.1. SUPPORTING IMMERSIVE/NATURAL INTERFACES	10	<b>4. THE POSITION OF EUROPE</b>	<b>37</b>
1.3.2. ENSURING SECURITY AND PRIVACY FOR PERSONAL DEVICES	10	<b>4.1. STRENGTHS</b>	<b>37</b>
<b>1.4. CROSS-CUTTING CHALLENGE 1: ENERGY EFFICIENCY</b>	<b>11</b>	4.1.1. STRONG EMBEDDED ECOSYSTEM	37
<b>1.5. CROSS-CUTTING CHALLENGE 2: SYSTEM COMPLEXITY</b>	<b>12</b>	4.1.2. PUBLIC FUNDING FOR R&D AND TECHNOLOGY TRANSFER	37
<b>1.6. CROSS-CUTTING CHALLENGE 3: DEPENDABILITY</b>	<b>13</b>	4.1.3. ONE OF THE BIGGEST MARKETS	38
<b>1.7. POLICY RECOMMENDATION: SUPPORTING INSTRUMENTS</b>	<b>14</b>	4.1.4. GOOD EDUCATION	38
1.7.1. SUPPORTING VIRTUAL VERTICALIZATION	14	<b>4.2. WEAKNESSES</b>	<b>38</b>
1.7.2. CONTINGENT FUNDING FOR COMMERCIAL DEVELOPMENT	14	4.2.1. EUROPE IS FULL OF HORIZONTAL SPECIALIZATION	38
1.7.3. CONTINUED FUNDING FOR ACADEMIC RESEARCH	14	4.2.2. LOSS OF COMPETITIVENESS IN SOME DOMAINS	39
1.7.4. INTERNATIONAL COLLABORATION	14	4.2.3. BORDERS AND DIFFERENT LANGUAGES	39
1.7.5. DEVELOPING PILOT-LINE FABRICATION CAPABILITIES	14	4.2.4. LACK OF VENTURE CAPITALISTS	39
1.7.6. LICENSING FOR COLLABORATIVE PROJECTS	14	4.2.5. WEAK ACADEMIA-INDUSTRY LINK	39
<b>2. MARKET TRENDS</b>	<b>17</b>	<b>4.3. OPPORTUNITIES</b>	<b>39</b>
<b>2.1. APPLICATION PULL: "THE INDUSTRIAL INTERNET"</b>	<b>17</b>	4.3.1. COST EFFECTIVE CUSTOMIZATION	39
2.1.1. "POST-PC" DEVICES: THE LINK BETWEEN HUMANS & CYBERSPACE	18	4.3.2. LEVERAGING FREE/CHEAP/OPEN INFRASTRUCTURE	39
2.1.2. NATURAL INTERFACES	18	4.3.3. SOCIETAL CHALLENGES	39
2.1.3. INTERACTION WITH THE PHYSICAL WORLD	20	4.3.4. CONVERGENCE	40
2.1.4. DATA DELUGE	20	4.3.5. MICRO- AND NANO-ELECTRONICS	40
2.1.5. INTELLIGENT PROCESSING	21	<b>4.4. THREATS</b>	<b>40</b>
2.1.6. PERSONALIZED SERVICES	22	4.4.1. COMPETING WITH FREE/CHEAP/OPEN INFRASTRUCTURE	40
2.1.7. SENSITIVE/CRITICAL DATA	23	4.4.2. FINANCIAL CRISIS	40
<b>2.2. BUSINESS TRENDS</b>	<b>23</b>	<b>5. CONCLUSION</b>	<b>41</b>
2.2.1. VERTICAL INTEGRATION WINS	23	<b>GLOSSARY AND ABBREVIATIONS</b>	<b>43</b>
2.2.2. ECONOMY OF SCALE FOR HARDWARE	24	<b>REFERENCES</b>	<b>45</b>
2.2.3. CUSTOMER LOCK-IN	24	<b>SUMMARY</b>	<b>46</b>
2.2.4. CUSTOMER/USER PROFILING	25		
2.2.5. FAB LABS	25		



# EXECUTIVE SUMMARY

Computer performance has increased by over 1,000-fold in the past three decades. This astonishing growth has fueled major innovations across all aspects of society. New advances in drug discovery and diagnosis, product design and manufacturing, transportation and energy, scientific and environmental modeling, social networking and entertainment, financial analysis, all depend on continued increases in computer system performance. Computing systems are so fundamental to today's society that they represent a basic resource, and form a strategic foundation for many of our most powerful and versatile tools and developments. Maintaining rapid growth in computing performance is key for tackling the societal challenges shaping Europe and assuring our global competitiveness in the future.

Yet today we are facing a new set of market trends and technological challenges to this progress, particularly within Europe:

- **Market:** The dominance of the computing systems market by desktops, laptops, and server PCs is waning and being replaced by a new market of smart *embedded* systems, *mobile* devices, and large-scale *data centers*. This new landscape is moving towards a *convergence* across embedded, mobile, and data center systems, where *global-scale applications* gather data from embedded systems and users, process it in large data centers, and provide customized, timely information to millions of users through their mobile devices or control our environment ("cyber-physical systems"). Addressing the *complexity* of system development and data processing at this scale will be critical for the next generation of systems and services. To support these applications, the market has seen a resurgence of *verticalization*, in which global companies strive to control the whole value chain from chip design through operating systems and applications all the way to end-user data services and sales. Cost and time to market are more than ever keys to market success.
- **Technology:** Energy has become the primary limiting factor in the development of all systems, whether due to the cost of energy or cooling in large systems or due to battery life in mobile devices. This has led to the rise of parallel and heterogeneous devices that trade off increased *complexity* and incompatibility with existing software for higher efficiency, and the appearance of "dark silicon", whereby portions of a device must be shut off to stay within the power limit. The necessity to develop energy aware devices and the ability to automate the optimization of applications for power efficiency has

become a necessity across all computing systems.

- **European Context:** Europe provides a strong embedded and low-power processor ecosystem, which will be critical for addressing power efficiency. This includes many companies in hardware and software development for both the industrial and commercial sectors. However, Europe also suffers from a high degree of horizontal specialization, which makes it difficult for companies to amortize the costs of development across the product chain.

This roadmap dives into the key computing systems challenges facing Europe in the next five years and provides recommendations for strategic research objectives for Horizon 2020. The roadmap leverages the broad academic and industrial expertise of the HiPEAC Network of Excellence to explore technology and market trends and identify promising directions for innovation. This document updates the 2011 roadmap with new trends in market verticalization, global-scale computing, and the impact of hardware design cost, while continuing to emphasize the difficulties of achieving energy-efficiency and programmability across devices from the cloud to mobile and embedded. From this analysis we identify *three strategic areas: embedded, mobile, and data center* and *three cross-cutting challenges: energy efficiency, system complexity, and dependability*.

## STRATEGIC AREAS FOR HORIZON 2020

### • Embedded systems

The traditional notion of an embedded system as a single-purpose device is rapidly changing as increased computing performance, connectivity and closer interactions with the world bring additional functionality and demands. To take advantage of this potential we need to rethink system architectures and programming models to optimize for energy, time constraints and safety and develop techniques to support portability of critical and mixed critical systems. Without such portability, the cost of certification for new computing platforms will prevent their uptake and limit our ability to leverage further advances.

### • Mobile systems

The shift from desktop PCs to mobile devices provides an incredible opportunity to rethink the human-computer interface and how we interact with technology. Innovations that provide more natural and immersive experiences will dramatically

improve the utility and productivity of mobile devices. Such developments require collaboration across all levels of the computing system: from human interaction to image processing and data mining, down to system architecture for efficiency and performance. Further, as mobile devices become increasingly integrated in our public and private lives, we will need stronger guarantees of privacy and security.

**• Data center computing**

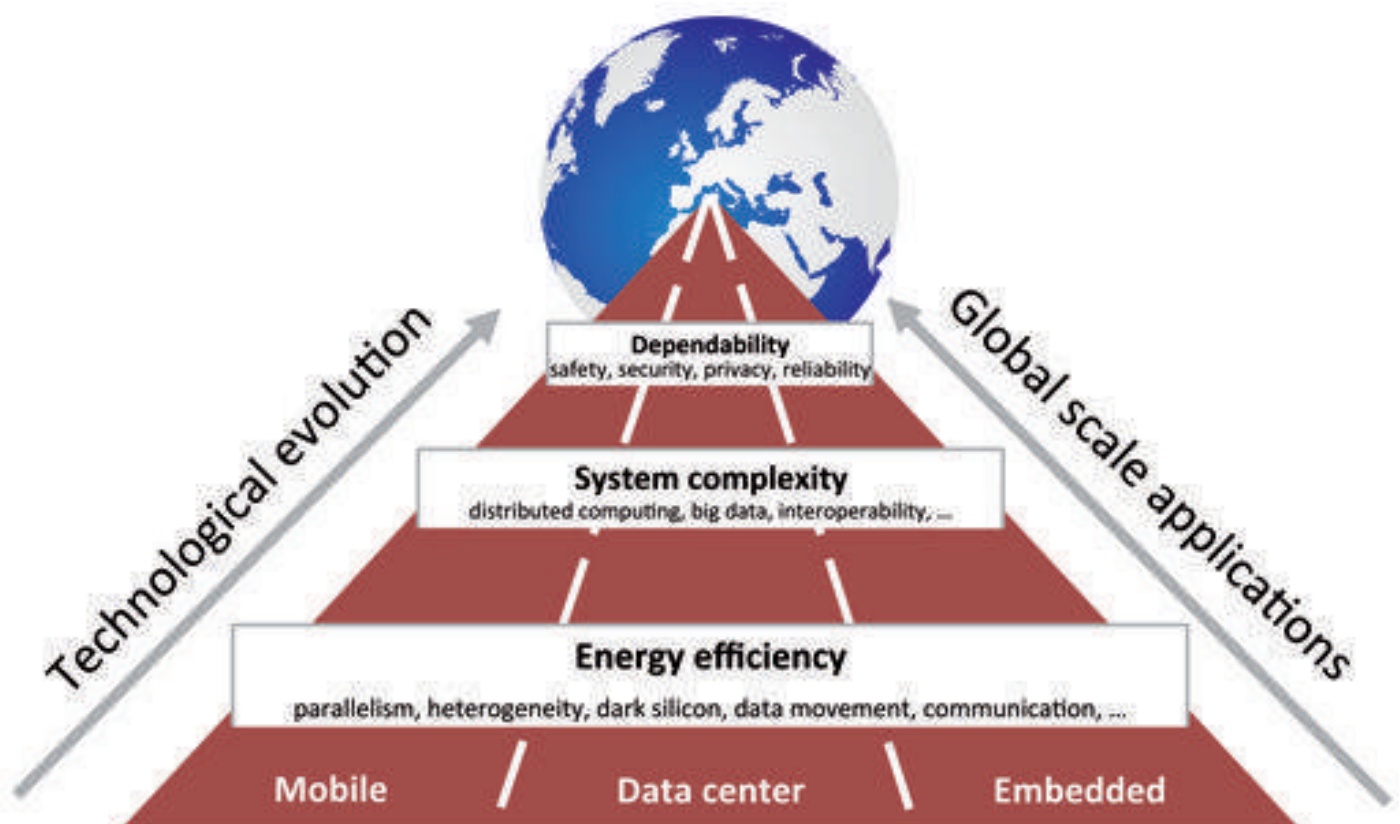
As applications become global-scale, Europe has an opportunity to lead the global market for data center technology. To be competitive we must develop the capabilities to process “big data” without increasing cost or energy. These challenges include architectures for handling massive unstructured data sets, low-power server modules and standards, network and storage systems, scalable software architectures, and micro servers. At the same time we must integrate these developments with techniques to ensure security, privacy, and compliance, while providing large-scale reliability, availability, and serviceability.

**CROSS-CUTTING CHALLENGES FOR HORIZON 2020**

**• Energy efficiency**

Systems today are limited in their performance by power used or dissipated. This has led to the combination of many specialized (heterogeneous) processors to increase efficiency. Unfortunately, this has also increased the complexity of programming to the point where it is prohibitively expensive for many applications. On top of this, the energy cost of moving data now exceeds that of computing results. To enable power efficient systems we must address the challenges of programming parallel heterogeneous processors and optimizing data movement, both for legacy applications and new computing modalities. We must also take advantage of the energy saving potential of new technologies, such as non-volatile memories, 2.5D and 3D integration techniques, new silicon technologies such as FinFETs and FDSOI, and new computing modalities such as stochastic and approximate systems and algorithms.

**The HIPEAC vision for Advanced Computing in Horizon 2020**



• **System complexity.**

Modern computing systems have grown to the scale where developers need to coordinate thousands of processors at once to accomplish complex tasks for large numbers of users and across massive data sets. To support this scale we need to develop tools and techniques to optimize for performance and ensure correct operation, while operating “at-scale”. On the hardware side, chips have become enormously more expensive to design, verify, and produce. Today’s cutting-edge technologies are so expensive that they are only affordable for devices that sell 10-100 million units. This cost limits product differentiation and makes market entry for new ideas extremely difficult. To overcome this we need to investigate new integration techniques that enable high levels of integration and differentiation without the cost of cutting-edge fabrication.

• **Dependability.**

Computing systems are involved in ever growing parts of our lives, from providing intelligent control for our transportation to keeping track of our friends and colleagues. As this involvement grows, we require higher levels of dependability. We expect computing systems to be trustable, reliable and secure from malicious attacks, to comply with all safety requirements, and to protect our privacy. Ensuring these properties will require more powerful methodologies and tools to design and implement dependable systems in a cost-effective way.

In addition to the six strategic areas and cross-cutting challenges listed above, HiPEAC recommends several supporting instruments

that will improve the effectiveness of Horizon 2020 in stimulating research, startups, industrial uptake. In particular, successful research results should be easily and directly transferred to European industry to address performance, cost, and time-to-market constraints. Practical recommendations include: the creation of “virtual verticals” by bringing together players across all parts of the product chain – including customers – to tackle key problems, not only scientific ones; encouraging the transfer to real marketable products through contingent funding for commercial development after the completion of STREP projects; extended-length funding of academic research to address the longer duration high-risk projects; carefully selected international collaboration outside of Europe to bring in complementary expertise; the development of pilot-line device fabrication capabilities accessible to SMEs and startups and bridging the gap between research prototypes and products; and analysis of open source licensing frameworks to ensure that EU-funded projects do not release IP under licenses that are incompatible with commercialization.

By addressing these research objectives through Horizon 2020 we will be able to ensure that Europe will continue to benefit from the promised growth of computing technology across all sectors. If we can take the lead in these areas, Europe will be in a position to develop new technologies and standards for the global computing systems market and transfer it to its industry. However, failure to address these challenges will significantly reduce our ability to leverage computing systems to improve our global competitiveness and tackle our social challenges.



The FP7 network of excellence in High Performance and Embedded Architecture and Compilation (HiPEAC) is Europe’s premier organization in the field of computing systems for coordinating research, improving mobility, and enhancing visibility. Created in 2004, HiPEAC gathers over 330 leading European academic and industrial computing system researchers from nearly 140 universities and 70 companies in one virtual center of excellence of 1200 researchers. HiPEAC encourages computing innovation in Europe by providing collaboration grants, internships, sabbaticals, and networking through the yearly HiPEAC conference, ACACES summer school, and the semi-annual computing systems week.

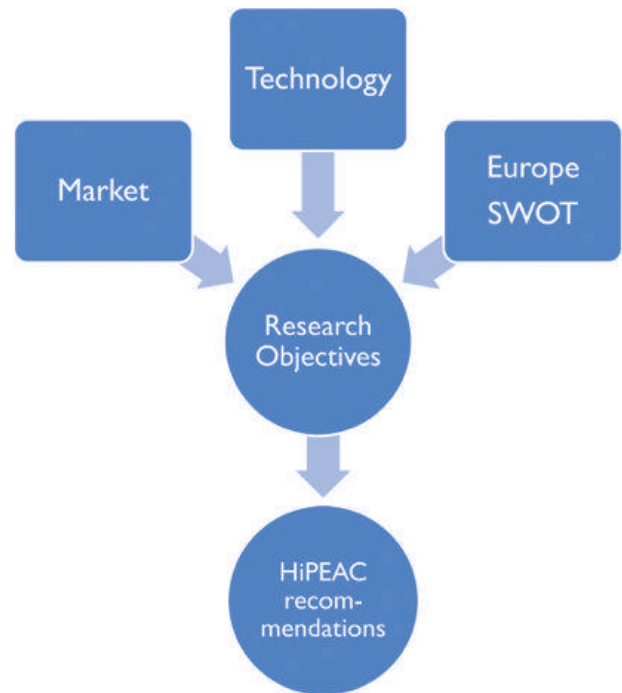
# PREFACE

Computing systems are universal. All aspects of public, private, and commercial life are affected both directly and indirectly by computing systems. Advances in computing are the key to the development of new domains and revolutionary technologies, such as personalized medicine, online social interaction, and immersive entertainment experiences. Across all of modern society, from manufacturing to agriculture, communications to energy, and social interaction to advanced science, computing systems are our primary tools for improving productivity, safety, well-being, and health. Computing systems play a major role in addressing our major societal challenges, including transportation, mobility, environment, energy conservation and education.

Yet, today computing systems are experiencing several dramatic shifts. Technological limitations are preventing the seemingly effortless performance increases of the past, while global-scale applications are placing computing systems into ever larger, more intensive, and more critical roles. At the same time applications and business trends are broadening the requirements for interoperability and flexibility. Devices have to integrate into global information processing chains, ranging from embedded systems interacting with the physical world, to mobile devices and data center servers for the cloud or High Performance Computing (HPC), always interconnected and communicating. The resulting dramatic increasing complexity, design and technology costs are limiting our ability to expand the impact of computing systems to new domains and applications.

This document is an update of the previous HiPEAC roadmap “Computing Systems: Research Challenges Ahead – the HiPEAC Vision 2011/2012”. It focuses on challenges for the coming five years, as formalized in concrete recommendations for the first years of HORIZON 2020. These recommendations are described in the first chapter of this document.

These recommendations are rooted and motivated in the subsequent parts: an analysis of the market trends, a discussion of the technology constraints and opportunities, and a review of Europe’s Strengths, Weaknesses, Opportunities, and Threats (SWOT) in the field of computing systems. From this information, the document identifies *three strategic areas: embedded, mobile, and data center* and *three cross-cutting challenges: energy efficiency, system complexity, and dependability*, which form the basis for the recommendations.



Structure of the document

## ACKNOWLEDGEMENTS

This document is based on the valuable inputs from the HiPEAC members, the teachers of the ACACES summer school 2012 and the participants of the roadmap meeting during the Computing System Week of October 2012. The editorial board, composed of Marc Duranton (CEA), David Black-Schaffer (Uppsala University), Koen de Bosschere (Ghent University) and Jonas Maebe (Ghent University), would like to particularly thank:

Angelos Bilas (Forth), Robin Bruce (Maxeler), Paolo Faraboschi (HP Labs), Grigori Fursin (INRIA), Christian Gamrat (CEA), Dimitris Gizopoulos (University of Athens), John Goodacre (ARM), Harm Munk (Astron), Marco Ottavi (University of Rome, Tor Vergata), Daniel Gracia Pérez (Thalès), Salvatore Pontarelli (University of Rome, Tor Vergata), Julien Ryckaert (imec), Per Stenström (Chalmers), Joseph van Vlijmen (ACE), Jennifer Sartor (Ghent University).



# ADVANCED COMPUTING SYSTEMS RECOMMENDATIONS FOR HORIZON 2020

Computing systems are experiencing a turbulent transition, full of both challenges and opportunities. The landscape of computing devices is changing: traditional desktop and servers are being replaced by mobile devices, smart embedded computing, and global data centers. Ubiquitous wireless communication provides constant access to unlimited amounts of data, thereby opening huge opportunities for novel services and applications, but also presenting significant risks to privacy and security. The convergence of embedded, mobile, and data center computing is moving us towards global-scale systems, which orchestrate thousands of machines, leading to a dramatic increase in system complexity.

In the future, the “computer” will be the global mesh of data, compute power, and interactivity stemming from the convergence of embedded computing, mobile computing, and data center computing. This integration can be referred to as the “*Industrial Internet*”, where the goal is to serve the user with trustworthy, reliable and safe services to support daily activities and needs.

At the same time that the world is shifting towards the convergence of embedded, mobile, and data center computing, the “free lunch” of continuous frequency scaling and power reduction of the past 50 years is over. The continuous performance improvements we have come to depend on have ceased due to the intrinsic power and frequency constraints of ever-smaller transistors. As a result, we are now faced with adapting to the complexities of heterogeneous parallelism in the quest for energy efficiency and increased performance. But the cost to leverage the advances of heterogeneous parallelism is too great for most applications due to the immaturity of software tools and techniques.

European industry and researchers are strong in several of the key technologies required for this future, in particular embedded systems and low-power processors. We believe that Europe should leverage these strengths through investment in tools and technologies to enable the future convergence of data center computing, mobile computing devices, and embedded devices and sensors.

To address these issues, HiPEAC has identified *three key strategic areas: mobile, embedded, and data center* and *three key cross-cutting challenges: energy, complexity, and dependability*. Addressing these strategic areas and tackling these challenges will require working across system and application boundaries, rethinking hardware and software interfaces, and investigating the impacts of technology and application evolution on algorithms and methodologies.

## STRATEGIC AREAS

### • Embedded computing

Embedded computing devices form the interface between the *physical world* and the *digital world*. They use sensors and actuators to measure and control physical phenomena, such as temperature, traffic, and electricity usage. Embedded computing devices have local computing power to preprocess raw data and extract salient features in order to reduce communication requirements or to locally process data in order to control actuators. Depending on the application, this computing power can be quite significant (e.g., intelligent surveillance cameras, Engine Control Units (ECUs) controlling the engine in a car) or negligible (e.g., thermostats). Embedded computing devices are ubiquitous. They are typically attached to the physical device they are monitoring and/or controlling, and communicate externally via a wired or wireless

network. Embedded computing devices are the building blocks of the *Internet of Things* and of *Cyber-Physical systems*.

#### • Mobile computing

Mobile computing devices such as smartphones and tablets form the interface between the *humans* and the *digital world* (or “cyberworld”). The role of mobile devices is to support natural interfaces, which provide humans access to the information stored in the digital world. Mobile devices are responsible for understanding their users’ environment, retrieving the most pertinent information, and visualizing the results in an intuitive manner. These devices will actively work to optimize their users’ experience by seeking out relevant information, providing a local sensing platform, offering new features.

#### • Data center computing

Data center computing is required to process the massive amounts of data generated by embedded and mobile computing systems, online transactions, and scientific simulations. These systems offer on-demand computing, storage, and scalability, and are backed by geographically distributed infrastructures and resources. The role of data center computing is to provide the computational capacity for analyzing and interpreting the data (e.g., real-time translation, forecasting, scientific computing, etc.) the storage capacity for recording it (e.g., customer activity, environmental conditions, experimental results, etc.) and ubiquitous access and scalability (e.g., to handle peak loads and provide reliability). Taken together, the many computing service providers and the network infrastructure constitute a “global data center”. Efficient communication networks are a mandatory element for off-loading computation and data to data center.

### CROSS-CUTTING CHALLENGES

#### • Energy efficiency

Systems today are limited in their performance by power used or dissipated. This has led to the combination of many specialized (heterogeneous) processors to increase efficiency. Unfortunately, this has also increased the complexity of programming to the point where it is prohibitively expensive for many applications. On top of this, the energy cost of moving data now exceeds that of computing results. To enable power efficient systems we must address the challenges of programming parallel heterogeneous processors and optimizing data movement, both for legacy applications and new computing modalities. We must also take advantage of the energy saving potential of new technologies, such as non-volatile memories, 2.5D and 3D integration techniques, new silicon technologies such as FinFETs and FDSOI, and new computing modalities such as stochastic and approximate systems and algorithms.

**System complexity.** Modern computing systems have grown to the scale where developers need to coordinate thousands of processors at once to accomplish complex tasks for large numbers

of users and across massive data sets. To support this scale we need to develop tools and techniques to optimize for performance and ensure correct operation, while operating “at-scale”. On the hardware side, chips have become enormously more expensive to design, verify, and produce. Today’s cutting-edge technologies are so expensive that they are only affordable for devices that sell 10-100 million units. This cost limits product differentiation and makes market entry for new ideas extremely difficult. To overcome this we need to investigate new integration techniques that enable high levels of integration and differentiation without the cost of cutting-edge fabrication.

#### • Dependability

Computing systems are involved in ever growing parts of our lives, from providing intelligent control for our transportation to keeping track of our friends and colleagues. As this involvement grows, we require higher levels of dependability. We expect computing systems to be trustable, reliable and secure from malicious attacks, to comply with all safety requirements, and to protect our privacy. Ensuring these properties will require more powerful methodologies and tools to design and implement cost-effective systems.

Future killer applications will come from the convergence of mobile and embedded interface devices and data center computing. This convergence will enable applications to dynamically redistribute computation and communications depending on the local environment and will operate “at-scale” to handle millions of global users and the processing of enormous data sets. For example, smart embedded computing in roads and vehicles will communicate traffic behavior and intentions. This data will be aggregated and processed in real-time by large-scale compute services. The results will then be personalized and displayed by in-car interfaces to optimize traffic flow. Similar scenarios exist for smart grids, smart cities, health monitoring, industrial automation, and national security. However, these systems must continuously adapt to local communications, power, and processing constraints, by adjusting how compute, sensing, and communications are distributed between the embedded systems, mobile systems, and global data centers. By developing the infrastructure and techniques to develop applications that span all three layers we will be able to benefit from the available data, interactivity, and compute power.

### 1.1. STRATEGIC AREA 1: EMBEDDED SYSTEMS

Europe is known for its large expertise and experience in the field of embedded systems and it should continue to be at the forefront of innovation in those domains. Competitiveness in this area is essential as the market for embedded systems dwarfs the mobile and data center markets, and is less dominated by non-European vertically integrated companies. Maintaining leadership in advanced embedded computing is a key opportunity for Europe.

Moving forward requires efficiently coping with timing constraints and safety, reducing certification costs, ensuring correctness in more complex systems, and supporting more complex devices. Emphasis should be given to how to design or leverage commodity computer system architectures for efficient support of mixed criticality and reduced certification costs.

### 1.1.1. COST-EFFECTIVE DESIGN OF EMBEDDED SYSTEMS

Designing embedded systems often requires hardware-software co-design for a product that will sell only tens of thousands of units. This small volume makes design too expensive with today's tools. To overcome this, more powerful design tools are needed that can reduce the time to adapt or create new designs, particularly through integration and cross-unit optimization.

The integration of multiple functions on generic hardware (often commodity multi-core platforms) is a compelling avenue for reducing design costs. However, this approach requires that the safety requirements of the individual functions can still be guaranteed in the presence of other active functions on the same hardware. There is a need for defining standard methodologies for the design of such mixed criticality systems and support for reducing their cost of certification.

Designing application-specific hardware and selecting key parts to integrate (e.g., processors, memory, analog, wireless, sensors, MEMS) can enable more optimized devices. However, as the development cost and mask cost for producing silicon chips is prohibitive for small product quantities, new techniques should be developed to enable developers differentiate hardware with a reduced development cost. Solutions such as integration of heterogeneous dies on a silicon interposer might allow differentiation while using the existing technological capabilities of the European semiconductor industry.

### 1.1.2. COST-EFFECTIVE CERTIFICATION

In sectors that deal with safety critical processes, embedded systems need to be certified before they can be deployed. The certification process is slow and costly. To increase the functionality and integration of embedded systems in safety critical applications we need to find ways to reduce the burden of certification. We should strive for a path to reduce certification, maintenance, and running costs.

For long-lived critical systems, such as airplane or train controllers, the cost of certification is so high that manufacturers are forced to stockpile replacement parts for the life of the product to avoid having to re-certify on new hardware. This requirement is extremely expensive, but could be avoided by enabling certification on virtual platforms. Such a certification would

allow manufacturers to substitute one certified virtual platform for another while maintaining certification. This would avoid the need to stockpile replacement parts and enable manufacturers to upgrade platforms to newer devices, which could reduce weight and energy requirements without requiring re-certification.

Similarly, the development of new computing paradigms and development methodologies that are driven by the constraints of safety-critical systems (e.g. taking into account time requirements, latency) can reduce the qualification and verification cost by ensuring correct-by-construction designs. This approach differs significantly from the majority of mainstream computing approaches, which strive for best-effort average case performance, and ignore many safety-critical functional requirements such as timing.

### 1.1.3. SECURE EMBEDDED SYSTEMS

The world of embedded systems has not been exposed to massive malware attacks in the past due to the general lack of external connectivity for embedded machines. However, today's systems are almost universally connected, and many high-profile hacking incidents have demonstrated that security needs to be improved in industrial and commercial embedded systems. Addressing this by hardware or software solutions will require a new focus on providing default security and updates for embedded systems ranging from industrial controllers to cars, pacemakers, televisions, refrigerators, smart meters, and even traffic lights.

## 1.2. STRATEGIC AREA 2: DATA CENTER COMPUTING

Applications are moving towards global-scale services, accessible across the world and on all devices. Low power processors, systems, and communications are key to computing at this scale. Europe's strengths in low-power embedded computing forms a strong foundation for becoming the leader in energy efficient data center computing infrastructure. This leadership will allow us to bring more European technologies and innovations into the global market for data center compute storage and cloud systems.

### 1.2.1. LOW-POWER MICRO SERVERS AND MODULES

Europe is currently the world leader for low power processors for mobile devices and is increasing its market share in microcontrollers for embedded systems and the Internet of Things. Moving beyond the embedded market, ARM has recently introduced its first 64-bit architecture to enable energy efficient server-class processors.

## 1.3. STRATEGIC AREA 3: MOBILE SYSTEMS

From a software point of view, efficient approaches for porting and optimizing existing data center software (typically written for x86 processors) to new energy efficient computing platforms will be required. Techniques for this transition include compiler optimization of binaries and libraries, dynamic recompilation, and even high-speed emulation or virtualization.

The goal of data center computing is to handle vast quantities of data efficiently and cheaply, storing and accessing data while minimizing cost and energy. This kind of workload is radically different from desktop, mobile, and HPC workloads. The hardware has to be adapted to the new workloads, and at the same time, the software has to be adapted to make optimal use of the resources offered by low-power server modules.

One opportunity to improve the hardware is the integration of efficient server-class processors with innovative technologies such as 3D integration (e.g. stacking memory chips directly on top of processors) or 2.5D integration (processors and memory integrated on a fast silicon interposer, potentially including optical interconnect). This presents an opportunity for Europe to produce very efficient computing “bricks”. Such “bricks” can form the basis for highly efficient cloud servers and low-power micro-servers. The micro-server form factor can provide extremely dense, small form-factor data center computing that can be installed anywhere without special infrastructure for cooling or power. These systems are of great interest for both large-scale data centers wishing to increase density and energy efficiency, and for small-scale operations, which may wish to keep their data physically on-site. However, the very high integration of these devices may require the development of new technologies for cooling (e.g., fluidic chip/wafer-level cooling), packaging, and interconnect.

### 1.2.2. NETWORK AND STORAGE I/O

Network communication and storage are at the heart of modern data centers and play a critical role in both scaling and efficiency. Achieving the required I/O rates to keep processors busy in large servers is essential for handling big data workloads. Achieving the required storage I/O performance, density, energy efficiency, and overall reliability, will require the use of new storage technologies, including non-volatile memories (NVM) and new approaches to scaling the I/O subsystem. Software and hardware techniques for coping with I/O latencies will also need to be improved. New technologies are needed for improving the performance and energy efficiency of high-speed networking, including host protocols, optical interconnects, network interfaces, links, switches, and topologies. In addition we must address the overhead of I/O virtualization in servers with increasing core counts.

The evolution of mobile devices will only accelerate as they become more powerful. New features will increasingly assist us with our daily tasks, anticipate our needs from our environment, and provide new capabilities, such as mobile payments and health and environmental monitoring. The value of these systems will come from their ability to integrate data from embedded sensors around the user and the world, process it in data centers, and deliver it directly to us in an intuitive and timely manner.

### 1.3.1. SUPPORTING IMMERSIVE/NATURAL INTERFACES

The key to the success of today’s mobile devices has been in their smooth integration of remote data with their users’ mobility. Natural interfaces, using human senses and interactions are important. Applications are beginning to use gesture, voice and image recognition and analysis as major sources of input to integrate more seamlessly with their user and environment. These more “natural” modalities augment information from traditional sensors such as Global Positioning Systems (GPS), compasses, gyroscopes, and accelerometers. On the output side, the graphics capabilities of mobile devices are making virtual or augmented reality practical. The end result of this progression will be devices that “disappear” into their environment by seamlessly adapting and providing the right data in the right context.

For the users, the mobile devices must support a smooth integration with data center processing and services to retrieve and process relevant information. However, much of the visualization and sensor pre-processing will happen on the mobile devices to reduce communication bandwidth. These functions, which deal with image processing, recognition and synthesis, will drive future hardware performance requirements, and delivering the seamless integration of embedded, mobile, and data center will demand significant advances in global-scale software development techniques.

### 1.3.2. ENSURING SECURITY AND PRIVACY FOR PERSONAL DEVICES

Ensuring the security and protection of personal data on mobile devices is a significant challenge. Separation of security domains between applications and between user contexts (e.g., professional and personal), private data, security against data theft, and secure payments are elements that should be developed. The solutions to these problems require cross-domain integration between the hardware, mobile device software, and data center applications. This will require the development of standards to enforce security and segregation of data, and provide for interoperability and validation.

## 1.4. CROSS-CUTTING CHALLENGE 1: ENERGY EFFICIENCY

All computing systems are power-constrained today: either by battery capacity, energy cost, or cooling. Addressing energy efficiency in embedded, mobile, and data center computing is essential for providing increases in performance.

Data center computing	Mobile computing	Embedded computing
Data centers/supercomputers are extremely power hungry. More power-efficient solutions (e.g. micro-servers) will reduce energy and cooling costs	Reduced energy consumption is needed to increase battery life for mobile devices	Energy usage determines the lifetime of battery-powered sensors and the extent of their computational capabilities
Power determines the maximum computing density in a data center and the cost of cooling infrastructure	Passive cooling limits the maximum power dissipation of a mobile device	Many sensors have limited ability to dissipate power (e.g., industrial or sealed sensors)

The power per transistor has ceased to decrease as we scale transistors to smaller and smaller sizes. This presents us with the difficult problem that modern chips can consume more power than we can either provide or afford to cool. To continue to take advantage of larger numbers of smaller transistors, we need to figure out how to build hardware and write software to make optimal use of these transistors.

### 1.4.1. OPTIMIZING DATA MOVEMENT AND COMMUNICATIONS

Today data movement uses more power than computation. We are now living in a world where *communicating and storing data is more expensive* (in both power and performance) *than computing* on it. To adapt to this change, we need to develop techniques for exposing data movement in applications and optimizing them at runtime and compile time and to investigate communication-optimized algorithms. Data movement needs to be viewed across the whole stack, not just within the processor. There is also a need to re-evaluate architecture decisions on memory and compute in light of new technology constraints and applications, for example, in light of 3D stacking, silicon photonics and new non-volatile memories.

### 1.4.2. PROGRAMMING HETEROGENEOUS PARALLEL PROCESSORS

The straightforward way to increase performance for a given power budget is to use special-purpose accelerators. Accelerators are processors optimized for particular tasks, which allows them to achieve significantly better energy efficiency at the cost of flexibility. However, a heterogeneous mix of different processors imposes a dramatic burden of complexity on the software, which now has to be optimized for different processors and has to intelligently decide where and how to execute. On top of this complexity, the variety and rapid rate of change of hardware make it very expensive to constantly adapt to new hardware generations. The current state of the art for heterogeneous development is not yet cost-effective.

To overcome the complexity and cost of heterogeneous system development, we need tools and techniques that provide power and performance portability, analyze software to provide actionable, high-level feedback to developers and runtime systems, and enable porting of legacy applications to modern hardware. The programmer should only express the concurrency of the application, and leave to the tools the mapping of the concurrency of the application into the parallelism of the hardware. Promising approaches include domain specific languages (DSLs), meta-languages, high-level intermediate representations with runtime optimization (JITs), combinations of runtime and static program analysis, and the pervasive use of metadata to provide context and hints to automatic optimization and mapping tools. Empirical program and architecture optimization combined with run-time adaptation and machine learning has demonstrated a high potential to improve performance, power consumption and other important metrics. However, these techniques are still far from widespread due to unbearably long exploration and training times, ever changing tools and interfaces, lack of a common methodology, and lack of unified mechanisms for knowledge building and exchange.

### 1.4.3. DEVELOPING NEW COMPUTING MODALITIES

Techniques such as biologically inspired approaches, approximate algorithms, stochastic and probabilistic computing, and in-memory computation have the potential to produce more energy efficient systems by relaxing accuracy requirements. Some of these technologies may also be particularly relevant for addressing silicon reliability issues in future process generations.

At the most basic level, many alternatives exist to current CMOS devices, but no clear winners have emerged as of yet [Berstein]. These alternative devices have the potential to reach switching energies three orders of magnitude lower than standard CMOS [Cavin]. Among these approaches, some may require a complete

rethinking of the whole design space. For example, field-coupled nano-computing [Imre06] offers ultra-low power computation and non-volatile behavior, but it uses the so-called “processing-in-wire” design paradigm, which is completely different from standard CMOS practices. The opportunity to rethink the whole system design in light of new technologies may lead to novel solutions for today’s problems.

## 1.5. CROSS-CUTTING CHALLENGE 2: SYSTEM COMPLEXITY

Future applications will require integration across large distributed systems, comprising thousands of embedded, mobile, and data center devices. Enabling future software and hardware developers and system designers to work quickly and effectively at this scale will require developing tools and techniques to efficiently optimize across system boundaries for performance and energy, ensure correctness and dependability, and enforce security and privacy requirements.

Data center computing	Mobile computing	Embedded computing
Scale-out to support global-scale applications and very large user bases	Optimize communications and computation for power consumption and user experience, particularly on wireless networks	For power efficiency massive parallelism is required (e.g., GPU, manycore)
Scale-out workloads and scale-out programming models. At-scale debugging and development, where developers support deployment across thousands of servers	Improved user experience through environmental awareness from sensor fusion and intelligent visualization	Reduced resource utilization but improved reliability and increased network integration. Certified programming languages and design paradigms
Intelligent processing, data mining, and complex analysis and computation	Media processing. Recognition, analysis and synthesis.	Feature extraction, hardware support for reducing power and real-time requirements
Large storage arrays containing unstructured data	Big media processing (mostly video)	Data capture and reduction

Programming the individual compute elements (data center, mobile, and embedded) in global-scale systems of thousands of nodes is simply not practical. We must develop methodologies and tools to allow the specification of the systems at a higher level, and relieve the developers of the need to manage the distribution of computation, communications, and storage. To enable cost-effective *global-scale applications* we need tools for simultaneous development and deployment “at scale”, that is, across thousands of machines at once. Developers need to be able to work “at scale”: debugging, testing, and profiling require languages, tools, and runtimes that support global-scale development and deployment. Tools for analyzing interdependencies and tracking performance through large multi-layered systems are mandatory. To be effective, these tools need to tie customer-focused performance metrics to software design (e.g., not just bits and bytes but interactions and product outcomes).

To accomplish this “at-scale”, we need to develop automatic or semi-automatic techniques to generate *system-level architectures*. These tools will use the application’s requirements (throughput, latency, processing, storage, etc.) and the technology constraints (data center capacity, compute and communications efficiency, energy costs, etc.) to select resources and connectivity that optimize for performance, complexity, and power, and map the applications to the required resources. At the system-level, this includes data center compute, storage, and communications capacity, partitioning applications across data center, mobile, and embedded compute elements, and optimizing tradeoffs in battery life, communications, and performance. Such design space exploration tools can reduce the time and cost of global-scale system development by allowing the developers to concentrate on functionality instead of the architecture required to achieve it.

**New methodologies** need to be developed to reduce the time it takes to efficiently implement applications and algorithms. Promising approaches include domain specific languages (DSLs), meta-languages, high-level intermediate representations with runtime optimization (JITs), combinations of runtime and static program analysis, and the pervasive use of metadata to provide context and hints to automatic optimization and mapping tools.

## 1.6. CROSS-CUTTING CHALLENGE 3: DEPENDABILITY

Global-scale systems are not only difficult to design and implement, but also difficult to prove safe and secure as are their components. However, given our increasing dependence on global-scale computing systems, we need to place more attention on the dependability of systems at all levels. A small mistake in a global system may easily lead to an incident with a worldwide impact.

Data center computing	Mobile computing	Embedded computing
Data security, privacy	Trustworthiness, security, anti-tampering	Physical security and safety, certification
Redundancy, statistical failure prediction, high availability	Replace on failure, ensure privacy of remaining data	Fault tolerance, error resilience
Mixture of batch and real-time (e.g., power grids, financial), latency sensitive for user interactivity	Low average latency for user interactivity	Closed loop sensors used in a real-time setting. Sensors in mobile devices often are not real-time

Dependability will play an increasingly important role in the next generation computing systems as they become more complex and interconnected, and deal with increasingly sensitive and critical information (e.g., health records, financial information, power regulation, transportation, etc.). In addition to the increase in dependability issues from a functional point of view, the electronic-level reliability of processors, memories, and communications channels is decreasing due to continued transistor scaling.

Privacy and security have become central issues as increasingly large amounts of data are stored in global data centers. Technologies and protocols must be developed at the hardware and software level to enable private data (e.g., medical, financial, industrial information) to be safely stored and processed in public and semi-public data centers. This should cover encryption and security for virtual machines at the hardware and software level, hardware and operating system support for data retention and deletion laws, and cryptographic support for law-enforcement and national security. Technologies should be designed with attention to legal compliance requirements.

For embedded systems, security has become as important as safety. Cars, planes and trains are being designed according to stringent safety requirements, yet insecure access to these

systems puts them at risk to hacking. High-profile international hacking of industrial controllers and demonstrations of remote hacks into cars, pacemakers, and hotel door locks have brought the issue of embedded system security to the forefront.

As the embedded systems interact with the physical world, they have to deal with non-functional requirements such as response time and reliability coming from various sources that have various level of trustworthiness. Methods to correctly fuse unreliable data from different sources in order to obtain a higher correctness at system level are important. These issues have rising priorities as we start to put significant degrees of high performance computing in real-time control loops (such as for smart grids) and their dependability needs to be ensured. Qualifying and ensuring correctness of a complex system composed of independently designed, distributed elements is a major challenge and new methodologies for proof, validation, and test are required.

For mobile systems, security is becoming more critical as people store more of their private lives and personal information on mobile devices. This leads to serious privacy and security problems when the mobile systems are hacked, stolen, lost, or broken. In all of these cases the data must remain secure to protect the individual. The amount of malware for mobile phones is currently growing exponentially as their popularity increases. Companies are having a hard time preventing employees from bringing their own mobile devices to work (BYOD, Bring Your Own Device), and the security implications of such untrusted devices on corporate networks is significant.

A global perspective on how to tackle dependability across all system layers and in all its facets is required. The implementation of dependability requirements should be done by automated tools based on the developer's specifications. Manual implementation of safety, security, and reliability features is error-prone and makes certification difficult. Retraining all software developers as security experts is simply not feasible and even more difficult than to retrain them as parallel programmers.

## 1.7. POLICY RECOMMENDATION: SUPPORTING INSTRUMENTS

HiPEAC's industrial partners have identified several ways in which publicly funded research partnerships are particularly relevant for companies:

- They enable investments in new or high-risk areas where the business case is not yet clear (this is particularly valuable for SMEs, which have limited financial resources to back exploratory product development and research).
- They provide resources to develop an ecosystem around existing technology, even if such an ecosystem will benefit competitors as well.

However, core technologies will continue to be developed internally due to intellectual property rights concerns.

The supporting instruments in Horizon 2020 should assist in areas where the markets will not, and should encourage academic-industrial partnerships, transfer of research results into pilot lines or to European companies ready to productize them, and cooperation across borders.

### 1.7.1. SUPPORTING VIRTUAL VERTICALIZATION

Supporting instruments should strive to bring together actors along the whole value chain to enable them to optimize across boundaries through close cooperation. Projects should explicitly seek players through most of the value chain from development to production and to the consumer in order to develop an ecosystem that can strengthen all companies together by delivering a compelling experience to the end user.

### 1.7.2. CONTINGENT FUNDING FOR COMMERCIAL DEVELOPMENT

The instruments should provide extended funding options contingent on a promising business case. This would allow projects to develop the project technologies and then re-evaluate whether they are commercializable. For example, by promoting the construction of demonstrators for field tests, or by a funding strategy with an option of 1-year commercialization extension for a subset of the partners.

### 1.7.3. CONTINUED FUNDING FOR ACADEMIC RESEARCH

The 3-year project is often too short for academic research. Providing a longer-term research instrument would allow academics to look into more risky projects. Currently the only option is

the European Research Council (ERC), which supports projects with a single principal investigator. The computing systems community in Europe is lacking the highly visible research labs that can attract international talent, serve as research hubs for industry, and provide incubators for new companies. There are a number of successful national initiatives like the excellence clusters funded by the German government. We believe that it makes sense to try to create and support a small number focused, physical 'eurolabs' with substantial funding for longer-term research.

### 1.7.4. INTERNATIONAL COLLABORATION

International collaboration outside of Europe is necessary when key competencies are not present in Europe. Such collaboration should be carefully chosen in order to maintain industrial competitiveness of Europe, while promoting access to external competencies and facilities.

### 1.7.5. DEVELOPING PILOT-LINE FABRICATION CAPABILITIES

Niche market (10k-100k devices) startups and SMEs have an extremely difficult time producing hardware or Application Specific Integrated Circuit (ASIC) to demonstrate their product. In addition to the cost of the Computer Aided Design (CAD) tools needed to design the chips, the low volumes make them unattractive customers for existing foundries. This effect stifles development of new hardware ideas. Creating a European "hardware competence center" that offers affordable access to design tools and silicon runs, could help the emergence of creativity and new markets for European companies. Such silicon runs could be done either as "test lots" to validate the design or on pilot lines specialized for small quantities (e.g. using direct etching techniques instead of masks, to reduce the major source of cost in advanced technologies).

### 1.7.6. SPECIAL LICENSES (NON-GPL) FOR COLLABORATIVE PROJECTS

Not all Open Source and Free Software licenses are compatible with all business models. In particular the GNU General Public License (by design) can make it hard for proprietary software vendors to leverage code. Releasing the results from publicly funded projects at least under an additional license that makes such integration easier, such as the Modified BSD or MIT license, can help alleviate these problems. Publically funded projects should also strive to maintain clear licensing and copyright guidelines and contact information for commercial use of their products.



# **THE HIPEAC VISION FOR ADVANCED COMPUTING IN HORIZON 2020: RATIONALE**





# MARKET TRENDS

Computing systems have become integrated global-scale systems with a tremendous economic impact. Today's driving applications require extensive integration of multiple systems and computing elements, with different levels of processing power, storage, and communication to provide rapid access to data and intelligent analysis. This section explores the market trends driving this development and their implications for computing systems in Europe.

The combination of massive amounts of unstructured data generated by numerous devices, a demand for intelligent processing and ubiquitous interconnection, strong guarantees of security, privacy, and trustworthiness, and an increasingly limited energy budget lead us to summarize the current trends in computing system applications as:

***“An environment of efficient, intelligent, trustworthy, interconnected devices at your service.”***

To meet the challenges posed by these trends we need to enable global optimization of storage, communications, and processing with orders of magnitude less energy than today.

## 2.1. APPLICATION PULL: “THE INDUSTRIAL INTERNET”

In the past three years we have witnessed a dramatic shift from “computers” to “mobile devices”. Smart phones and tablets now

dominate the consumer market, while desktop computer sales are decreasing, leading to a “Post-PC world” with the key characteristic of *pervasive connectivity*.

The key categories of computing devices in this new world of pervasive connectivity are:

- **Data center computing:** clouds, servers, HPC systems, and large-scale storage systems, which are located in large-scale data centers and accessed remotely by users. The trend is towards providing computation as a utility that can be purchased from multiple vendors to fit the needs of the user.
- **Mobile computing:** PCs, tablets, mobile devices, and smart phones form the primary means of human-computer interaction. The trend is towards increasingly “humanized” interfaces, such as the use of natural language communication, vision, and gesture recognition.
- **Embedded computing:** “invisible” (for users) computing elements and (deeply-) embedded systems, including low-power, often single-purpose, computing systems that are typically tightly coupled to sensors and actuators. The trend is towards higher levels of functional integration while maintaining critical performance.

These three types of computing devices combine to form the *Industrial Internet*: “a network that binds together intelligent machines, software analytics and people... [to improve] productivity and growth.” [Annun12] The integration of these devices requires pervasive connectivity for nearly all tasks. For example, a mobile phone may read a Radio-Frequency Identification (RFID) tag, look up the product information from a server, and use that information to search for product reviews for the user. Without all three types of computing and ubiquitous connectivity, these services cannot function.

## THE “INDUSTRIAL INTERNET” REVOLUTION

According to the CEO of General Electric, “*The real opportunity for change is still ahead of us, surpassing the magnitude of the development and adoption of the consumer Internet. It is what we call the “Industrial Internet,” an open, global network that connects people, data and machines. The Industrial Internet is aimed at advancing the critical industries that power, move and treat the world.*”

“*The vast physical world of machines, facilities, fleets and networks can more deeply merge with the connectivity, big data and analytics of the digital world. This is what the Industrial Internet Revolution is all about.*”

“*The Industrial Internet leverages the power of the cloud to connect machines embedded with sensors and sophisticated software to other machines (and to us) so we can extract data, make sense of it and find meaning where it did not exist before. Machines – from jet engines to gas turbines to CT scanners – will have the analytical intelligence to self-diagnose and self-correct. They will be able to deliver the right information to the right people, all in real time. When machines can sense conditions and communicate, they become instruments of understanding. They create knowledge from which we can act quickly, saving money and producing better outcomes.*”

From <http://gigaom.com/2012/11/28/the-future-of-the-internet-is-intelligent-machines/>

### 2.1.1. “POST-PC” DEVICES: THE LINK BETWEEN HUMANS AND CYBERSPACE

Tablets and mobile devices are gaining momentum and are quickly surpassing PCs as the primary form of consumer computing. The combination of their convenience, intuitive interface, and ability to access large amounts of online data has led to the Post-PC revolution. Consumers no longer perceive their smart phones or tablets as “computers”, but rather as information and entertainment devices.

“*Look at PCs...we sold more iPads last year than market leader HP sold of their entire PC lineup. There were about 120 million sold last year...projections suggest 375 million within four years. Will exceed PC sales.*”

Tim Cook, Apple CEO

Mobile devices are now the main interface between humans, the real physical world and the online “cyberspace”. Pictures, movies, location data, etc., are created, manipulated and viewed through these devices, and count for a major part of the data exchanged on the internet. One important factor of their success is that they are “*bigger on the inside than on the outside*” thanks to their constant connectivity. This allows the user to have access to all the data on the internet, rather than being limited to the data stored locally. The success of mobile devices also resides in the fact that their interface is more natural for humans than a keyboard and mouse. Touch screens, gesture recognition (e.g. Microsoft’s Kinect), voice control (e.g. by Apple’s “Siri”), haptic feedback, and life-like imaging are providing natural interaction and are following the growing expectation that technology should be humanized and natural. We can expect that vision will be the next step in that direction, with devices “understanding” image content and their visual surroundings, and with new display technologies, glasses, and active-contacts leading to compelling reality augmentation.

Next generation devices will use more and more natural sensory interaction to improve their understanding of our intentions and requirements, delivering a personalized and unique interaction. These devices will be aware of their surroundings, which will enable them to empathize with and understand their user. This awareness will require merging and understanding data from an increasing number of other devices and sensors, including the user’s personal data, to provide the most compelling user experience. Such trends will drive advances in processing power and data security and privacy.

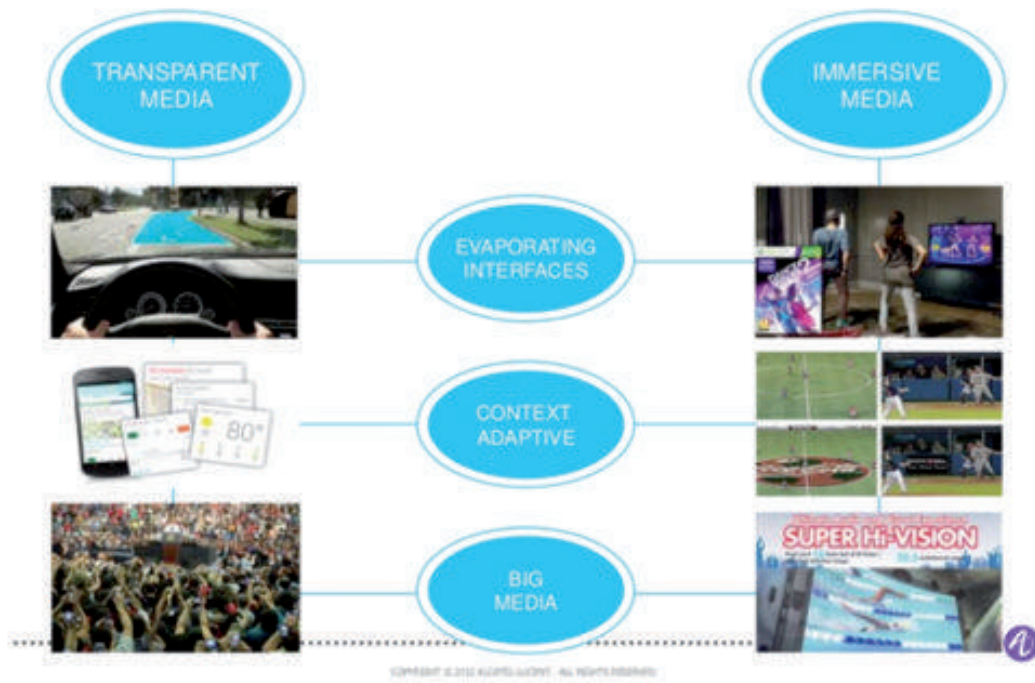
### 2.1.2. NATURAL INTERFACES

Computing interfaces are becoming more natural. Users increasingly expect to use voice and touch-based controls with mobile phones and tablets. Text-to-speech and voice recognition services are improving rapidly and are being applied everywhere, from the automotive industry to public transportation, and as automated real-time translation services.

In the future, user interfaces will blend into the environment to form so-called “evaporating interfaces”. Touch screens will have haptic feedback to simulate textures. Devices will be able to analyze and understand their auditory and visual environment, and sample the air to detect pollen and pollutants. Transparent heads-up displays are already used to overlay augmented reality in planes and high-end cars. Realistic avatars with advanced artificial intelligence will better understand the user and assist more efficiently – making manuals superfluous and man-machine interfaces more natural. Vision and image understanding will undoubtedly be the next major direction for device humanization and the entertainment industry will probably lead this evolution.

# FUTURE MULTIMEDIA EXPERIENCES

## KEY TRENDS



Ingrid Van De Voorde, Multimedia Technologies Research Domain Leader, Alcatel-Lucent, Antwerp iminds 2012

## Global integration of services



Cyber Systems (cloud, HPC, data servers), mobile and Cyber-physical Systems are all interconnected  
 Courtesy Jan M. Rabaey, UC Berkeley, updated for this HIPEAC vision

### 2.1.3. INTERACTION WITH THE PHYSICAL WORLD

Less apparent to the public is the explosion of embedded systems. These “small” systems are everywhere to gather information from the physical world, and yet they are nearly invisible. From the temperature sensor connected to a home thermostat to the dozens of embedded processors in cars, embedded data processing systems are invading our life by the billions. Their processing and communication capabilities are increasing by linking them to each other and to the Internet. This evolution enables increased functionality, such as internet-connected home thermostats that automatically learn your temperature preferences and allow remote operation. Indeed, it is becoming increasingly difficult to identify a single device where a complete application is running. Processing is now generally distributed between local pre-processing, processing for a user interface on a mobile device, and a live connection with the Internet for further processing or data collection from distant servers.

Embedded devices are increasingly interconnected, and communicate without human intervention. Kevin Ashton expressed this trend as follows: *“If we had computers that knew everything there was to know about things – using data they gathered without any help from us – we would be able to track and count everything, and greatly reduce waste, loss and cost. We would know when things needed replacing, repairing or recalling, and whether they were fresh or past their best. The Internet of Things has the potential to change the world, just as the Internet did. Maybe even more so.”*

This global integration offers services and capabilities that were not possible before. For example, it can be used to save energy and resources (“smart grids”) at a global or local level. In the case of electric cars, energy efficiency can be realized by creating an ecosystem where the car can be either recharged from a household power supply or provide energy to the local grid, depending on the current energy price. To optimize efficiency, the car uses information on traffic conditions and weather collected by embedded sensors to plan its route.

Computing systems will play a major role in healthcare, by monitoring wellness, identifying hazardous situations, and calling for help in emergencies.

More human-centric embedded systems will continue to evolve as well. In the near future it will be possible to check vital health parameters in real time, monitor activities and administer drugs in an adaptive manner. The performance of the human body may even be improved by enhancing certain functions such as hearing, vision and even smell (e.g. by warning a person to the presence of toxic particles). The augmentation of the human body through intelligent technology is colloquially referred to as “human++”.

### 2.1.4. DATA DELUGE

For the past decade, society has been generating exponentially increasing amounts of data. Commercial data sources include financial transactions, search histories, and product information. Public agencies contribute medical records, population databases, and legal and legislative data. Science and engineering routinely undertake large-scale simulations for weather prediction, drug discovery, product development, and record raw data from advanced sensors and experiments. There is also an exponential growth of data coming from data sensors (results of seismic data for oil exploration, genomic data in biology, etc.). Individuals are publishing an unprecedented amount of personal data, greatly encouraged by the explosion of social networking. This massive growth of data is forcing us to reevaluate how we work with data in computer systems.

#### THE DATA CREATION BY RADIO ASTRONOMICAL FACILITIES

The LOw Frequency ARray (LOFAR), located in the Netherlands (Astron facilities), is a phased array radio telescope, in which the signals from a large collection of low cost antennae are combined into images. LOFAR has ~45,000 antennae, and produces 200-300 TB of data during a 6 hour observation (Throughput of 78 Gbit/s). The Square Kilometer Array (SKA) telescope, to be built in South Africa and Australia, will add at least one order of magnitude to these numbers.

This scale of scientific data is causing significant challenges for the radio astronomy community:

- Shipping high volumes of data (petabytes/day): the raw observation data is pre-processed locally to reduce the volume, but needs to be shipped to different centers for further analysis.
- Locality of processing vs. distributed processing: how to determine criteria (cost, energy, bandwidth)
- Complexity of process management: many antennae have to be controlled. Many data streams have to be routed to the processing facilities. For example, during pre-processing, those streams have to be coordinated in real time (the networks have insufficient buffer capacity to absorb more than a few “minutes” worth of data.)
- Data processing: the amount of computation per output is very large. There is an intricate interplay between hardware architecture and algorithms to process the large amounts of data in a “reasonable” amount of time,
- Graceful degradation/robustness: how to deal with failing hardware?

This data is not only a byproduct of our digital society, but also a valuable resource. Buried within this data are key insights into human behavior, market trends, diseases, engineering safety, environmental change, and the basic workings of physics. Yet with such a massive deluge of data it is becoming increasingly difficult to get the most from it, and to do so in a timely manner.

*“The new resources of humanity are the data.”*

To exploit this data we need to be able to analyze it and respond in real time. Rapid analysis is key not only for financial trading (where a microsecond advantage can make the difference between a profit and a loss) and safety systems (where unknown delays can cause disaster), but also for delivering compelling customer experiences (by identifying trends early) and for society as a whole (identifying pandemics and optimizing transportation). Other opportunities are also emerging: massively instrumenting the energy grid in order to detect problems early to avoid general disasters. Reliable response time concerns are growing in data centers and HPC to avoid cascading delays caused by unpredictable response times.

The challenge for the next decade will be in coping with this humongous increase in data and the simultaneous demand for processing it faster and in due time.

### DATA DELUGE: A FEW FIGURES

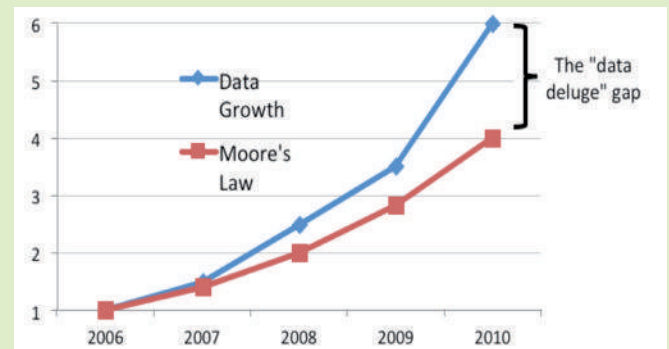
The term “Data Deluge” was coined in 2003 [Hey03] in the context of scientific data management to describe the massive growth in the data volume generated in research (and by scientific instruments), which was rapidly dwarfing all the data previously collected in the history of research. Since then, the ability to generate vast quantities of data has outpaced the infrastructure and support tools. This is true for scientific data, digital media (audio and video), commercial transactions, social networks, legal and medical records, digital libraries, and so on. The need to analyze, organize, and sustain “big data” is one of the highest priorities in information technology across disciplines, organizations and geographies [E10].

In 2010 the world generated over 1.2 Zettabytes ( $10^{21}$  bytes) of new data, 50% more than it had in all of human history before that. To put this in perspective, 120 Terabytes of new data was generated in the time it took to read the previous sentence. For example, Microsoft Update and Windows Update push out a Petabyte of updates monthly. Cisco predicts that by 2013 annual Internet traffic flowing will reach 667 Exabytes. A social network like Facebook produces 10TB/day of data, with Twitter is not far behind (7 TB/day); each of the 4.6B mobile phones

and 30B RFID tags produce several events per seconds that need to be stored, processed and analyzed. Likewise, the 2B Internet users also generate a variety of events that can have important value in areas like statistics, demographics or marketing. And the 50B connected devices expected by the year 2020 will cause all of the previously mentioned figures to balloon even further. Domains like gaming and other virtual worlds or augmented reality are also turning into massive data management problems.

In the scientific community, discovery has turned into a data-driven process, which represents a relatively new fourth paradigm in science [Hey09], next to the empirical, theoretical and computational models. The problem is that with current increases in computation, it may take over a decade to gain some understanding of what has already been archived from the most important scientific experiments. All fields of science (astronomy, physics, energy, medicine, drug discovery, climate, public health, etc.) are completely swamped with data, which requires major breakthroughs in repositories, storage, and computing architectures, all of which are central to the HiPEAC mission.

If we compare this growth with Moore’s law (transistor density doubling every two years, see below), it is clear that data is on a higher exponential growth than computation capacity, and this unprecedented trend is forcing us to reevaluate how we work with data in computer systems.



*Data growth vs. Moore's Law trends in the last 5 years. Data “deluge” means that we are heading towards a world where we will have more data available than we can process.*

### 2.1.5. INTELLIGENT PROCESSING

As we populate the physical world with devices and sensors that are connected to the cloud, they generate a constantly growing stream of unstructured data, audio, video, and physical measurements. To make the most of this information, we need systems that can analyze the data to extract useful information, as well as detect threats (for security – like detecting intrusion or abnormal behavior) and faults (out of tolerance in factory automation applications, for example). As the volume of data is growing exponentially, we need



From Dr. Rama Shukla, VP Intel Corp. *The New Age of Computing Continuum Experience*

to be extremely efficient in our data processing to keep up. One major example of this is video processing: robustly analyzing video streams is a very important technical challenge for computing systems. Currently, millions of surveillance cameras are producing footage that is never analyzed because there is no economical way of doing so. Human observers are too expensive and the current analysis programs are not reliable enough. Video analysis currently works well only for specialized tasks in reasonably structured environments (recognizing traffic signs, license plates, dedicated inspection tasks in factories, etc.). Robust video analysis is likely to become one of the next killer applications of computing systems. It is key enabling technology for autonomous cars, reliable inspection systems, security, patient monitoring, etc. The development of robust video analysis and its pervasive application would lead to huge economic savings, for example by enabling productive commutes via autonomous cars. Similar advances are needed in other domains such as speech recognition, customer analytics, and sensor data processing. On the personal side, automated indexing of personal pictures or movies is just starting with basic face recognition. This technology will enable easy retrieval by content and context, without the need to specify name and date.

The intelligent processing of data will be distributed across embedded sensors, mobile devices, and data center computers. This is required to minimize communication (e.g., pre-processing video to extract salient features) and to provide greater processing power when needed (e.g., for large-scale data analyses from many sensors). Such distributed processing may make sense for

privacy and security as well. For example, a camera that monitors if an elderly person has fallen in his or her bathroom should contain enough local processing to detect a fall and transmit an alarm, without having to send video images even outside of the camera. On other systems, processing at different levels in the device minimizes data transmission to the relevant and useful data and eases the fusion of data from different sources in order to generate the most appropriate result. This distributed and intelligent processing represents the new frontier that our computing systems will have to deal with in the years to come.

### 2.1.6. PERSONALIZED SERVICES

Services are becoming increasingly personalized, both in our private and professional lives. Our preferences are taken into account when accessing web-based services and are taken into account when renting cars or hotel rooms. Other examples are personalized traffic advice, search engines that take our preferences and geographical location into account (with the downside of creating so-called “search filter bubbles”, which means that the user only gets to see what he is interested in, leading to a distorted view of reality), music and video sources presenting items fitting our personal tastes and in the format that best suits our mobile video device, and usability adaptations for disabled people.

Personalized video content distribution is another case of ever-increasing importance. Video streams can be adapted to the viewer's point of view, to his or her personal taste, to a custom angle in case of a multi-camera recording, to the viewer's location, to the image quality of the display, or to the consumer profile with respect to the advertisements shown around a sports field. Another growing area is personalized medicine in which doctors try to optimally match disease types (e.g. hard to cure cancer types) and advanced drugs based on large databases of patient records. The common challenge for all personalized services is that they need lots of compute power because the content needs to be adapted per (type of) user.

### 2.1.7. SENSITIVE/CRITICAL DATA

Computing systems today contain nearly all details of our personal lives: from financial transactions, to eating habits, location information, medical records, and personal communications. This information can be used to infer what you do and like and with whom you interact. The potential privacy violations from having so much private data online are enormous. Large corporations are already aware of many of these risks and attempt to protect themselves with firewalls and private networks, preventing use of personal devices for professional activities, and blocking access to online data sharing services (like Gmail, Dropbox, Facebook, etc). The average user is not too concerned and only realizes that there is a problem when his or her identity is stolen and used for financial transactions, or when a smart phone goes missing or a virus takes control of his or her computer. Yet despite this blissful ignorance, companies are paying for the services they provide to the consumer by owning the personal profile information and selling it to the highest bidder in the form of advertising. The use and protection of such sensitive private data is neither well-understood nor regulated.

For computing systems controlling safety-critical applications such as power plants, cars, trains and planes, the system also determines our personal safety. An attack on a safety-critical system can easily lead to dangerous situations. The effects of this have been seen with military cyber-attacks in the Middle East and numerous security flaws have been found in the industrial controllers used for much of our basic infrastructure. Demonstrations of the ability to hack pacemakers and stop patient's hearts and the routine announcements of personal data theft from private companies show the direct personal impact poor cyber-security can produce [*Pacemaker*]. The potential for such damage will only continue to increase as we develop an increasingly integrated and automated society.

From Dr. Rama Shukla, VP Intel Corp. *The New Age of Computing Continuum Experience*

## 2.2. BUSINESS TRENDS

### 2.2.1. VERTICAL INTEGRATION WINS

ICT companies that have gained control over the complete value chain from hardware to consumer have shown tremendous growth recently. Apple has shown the most dramatic success by developing both hardware and software and controlling the content and a significant part of the retail distribution channel. Other companies are trying to implement similar models: Microsoft is developing retail shops, an online App Store, and increasingly sells hardware. Similar stories apply to Google, Samsung and Amazon. Europe has, relatively speaking, few such companies, with most of them actually having become less vertical over the last decade (Philips, Siemens) through spin-offs. In the war for patents and talent, major companies buy niche players mainly for their patent portfolios. Unfortunately, Europe is full of horizontal specialization, which makes it difficult to compete with massive vertically integrated companies. Hundreds of innovative SMEs are bought by large (mostly non-EU) companies every year.

An added advantage of vertically integrated companies is that global optimization – one of the technical challenges of this roadmap – is far easier for vertically integrated companies because they control a much larger portion of the system and can avoid the overhead of standard interfaces and unneeded abstraction layers. While standard interfaces and abstractions are needed to enable third party integration, they necessarily lead to more costly solutions in terms of integration costs and power and performance. Controlling the complete chain allows global optimization, at the expense of some choices that would not be optimal if considered independently. For example, Apple products require less memory than Android's because the software can be tuned for the particular hardware, and the hardware can be tuned for the operating system features.

*“For many year, this idea of “vertical integration” was out of favor. People thought it was kind of crazy, but we never did and we continued to build. This is something you work decades for. I think there are people trying desperately to catch up, and they’re finding it very difficult to do.”*

Tim Cook, Apple CEO

Furthermore, in some domains, such as development tools, it is extremely hard to compete with open source alternatives that may be less sophisticated, but free and good enough to do the job. Vertically-integrated companies can afford to subsidize the development of tools through the sales of products and services. The tools are a means for selling more products, not a product that needs to pay for itself. But subsidizing the development of advanced tools is expensive, and can only be done if the company is very large. Examples of this today include the LLVM compiler from Apple, the Go language from Google, and various performance tools released by companies such as Facebook and Amazon.

Vertically-integrated companies may also be more resilient to economic downturns because they do not have to move or follow a separate ecosystem. Due to their verticality, one weak segment can be compensated by other parts of the chain. This allows them to more easily tolerate market cycles, such as variations in prices of Flash memory or DRAM. Further, their large volume allows them to buy components at significantly reduced prices compared to less vertically-integrated competitors.

### 2.2.2. ECONOMY OF SCALE FOR HARDWARE

The cost of silicon development in advanced technology incurs non-recurring engineering (NRE) cost that can only be amortized when sold by millions. This implies that companies will no longer be able to afford to develop custom chips for any but the largest volume markets. General-purpose processors and FPGAs (Field-Programmable Gate Arrays) will not be affected by this constraint because of their flexibility, but ASICs (Application Specific Integrated Circuits) will be. Therefore, more and more systems will be forced to reuse off-the-shelf components instead of custom chips. These devices will have to differentiate themselves via software or by using different combinations of existing chips. This is a limiting factor for small companies that require the functionality or performance of a specialized chip in their product.

In a similar domain, this economy of scale is one of the successes of Apple: instead of having hundreds of products, they only have a handful of products (typically one in the mobile domain, with the “lower end” being the previous generation). This drastically reduces the development cost and amortizes the development cost by the millions of pieces sold. The hardware is then “personalized” by the software. The hardware is a white (in fact

black) sheet (slate) sold identically in millions, with no specific buttons. The user customizes the hardware by choosing its functionality via software, where each application can have its own user interface that the user chooses (that is why there are hundreds of Apps doing nearly exactly the same function, only differentiated by different user interfaces, depending on the taste of the user, who finally has a *personalized* device from impersonalized hardware).

### THE EXPLODING DESIGN COST

- Total SoC design costs increased 39% from the 32nm node to the 28nm node and are expected to increase 29% again at the 22nm node.
- Total software design costs increased 42.5% at the 28nm node and are forecast to show a CAGR (Compound annual growth rate) of 138.9% through the 14nm node.
- Derivative SoC silicon designs allow designers to accomplish their solutions at a fraction of the cost compared to first time efforts at the same process node when it first becomes commercially available.
- Costs for an Advanced Performance Multi-core SoC design, continuously done at the 45nm node, will experience a negative CAGR of 12.5% by the time the 14nm process geometry becomes commercially available, showing that subsequent designs at the same node become less expensive over time.
- 28nm silicon with a \$20 average selling price is required to ship 6.521M units to reach the breakeven point.

(From <http://www.semico.com/press/press.asp?id=299>)

### 2.2.3. CUSTOMER LOCK-IN

Companies try to lock their most profitable customers into their ecosystem: if you already have or use their products, they try to convince you to buy more products that belong to the same ecosystem. Notorious examples are providers that offer quadruple play deals: cable television, telephony, mobile, and Internet. Specialized service providers can offer cheaper deals for one service, but have a hard time competing against quadruple play packages.

Other examples are the vertical ecosystems by Apple, Microsoft and Google. Using a cell phone, tablet computer and computer of a single ecosystem is not necessarily cheaper, but gives the benefit of seamless interoperability. Since people seldom replace all of their devices at once, there must be a strong incentive for a user to go through the pain of moving from one ecosystem to another, one device at a time. It is clear that these companies also try to extend their ecosystem beyond the traditional computing devices, with television and cars being obvious targets for the coming years.

Another form of ecosystem lock-in is that social media websites attract customers with personalized services (Facebook, LinkedIn, Google+,...) that make it quite difficult to move out of their system by not allowing easy deletion of accounts, by making it cumbersome to transfer information from one website to another, etc.

#### 2.2.4. CUSTOMER/USER PROFILING

All businesses have discovered the value of gathering information about their customers/users, either for targeting advertisements or sales. Companies do their best to encourage the customer to provide even more personal data. A worrisome trend is that many websites today refuse to do business unless the consumer provides information unrelated to the core business transaction. Further, the consumer has little control over how all this information is used/exploited afterwards. The extreme form of this trend can be found in the social media where users voluntarily give private information to companies like Facebook, Google and LinkedIn. They even accept that these companies use this information to make money by sending them focused advertising. It has been demonstrated repeatedly that this information sharing leads to highly unwanted effects, but nevertheless millions of people do not seem to care.

#### 2.2.5. FAB LABS

Another potentially interesting trend is the opposite of producing mass market goods done by giant companies: it is, on the contrary, to make yourself (or somebody who has the skills) produce customized goods in small quantities, a sort of comeback of the small business sector and craft industries. This rebirth might be possible by “fab labs”.

##### THE EMERGING FAB LABS

*Fab Labs = fabrication laboratories are small-scale workshops with modern computer-controlled equipment that aim to develop Personal Fabricators (Fabbers).*

From Stefan Merten:

*Fabbers are universal fabrication machines which materialize three-dimensional things from digital data, typically by “baking” some amorphous, fine-grained material to allow for the constructing of things which are difficult to create otherwise.*

*Fabbers have some interesting features:*

- Link digital data and material world closely: Bringing the logic of digital data to the material world
- Universal for (parts of) material production: Like computers for information: “One machine to rule them all”
- Allow material production for individual needs, hence the concept of Personal Fabricators

(From [http://p2pfoundation.net/Fab\\_Labs](http://p2pfoundation.net/Fab_Labs) )

The fab lab is the answer to the willingness of people to have personalized devices and also their creativity. 3D printing capabilities are being investigated in medicine to “print” organs [TimeTech] or even to build a lunar base [ESA].

In order to work, the fab lab will require a lot of compute power, not only to model and simulate objects, but also to control the tools like 3D printers. Modelers and 3D tools will require high interactivity that can be offered by local high performance computers or remotely on the cloud if the latency is acceptable. This is a drive for micro-servers, which provide significant computational power at an affordable cost, and programming tools for low-cost high-speed embedded processors.

##### 3D PRINTERS: THE DIGITAL FABRICATION REVOLUTION

Neil Gershenfeld, Director of MIT’s Center for Bits and Atoms, believes that personal fabricators will allow us to do just that and revolutionize our world.

He explores *“the ability to design and produce your own products, in your own home, with a machine that combines consumer electronics with industrial tools. Such machines, Personal fabricators, offer the promise of making almost anything - including new personal fabricators - and as a result revolutionize the world just as personal computers did a generation ago.”*

*“A new digital revolution is coming, this time in fabrication. It draws on the same insights that led to the earlier digitizations of communication and computation, but now what is being programmed is the physical world rather than the virtual one. Digital fabrication will allow individuals to design and produce tangible objects on demand, wherever and whenever they need them. Widespread access to these technologies will challenge traditional models of business, aid, and education.”*

(From Neil Gershenfeld, <http://www.foreignaffairs.com/articles/138154/neil-gershenfeld/how-to-make-almost-anything>)





# TECHNOLOGY CONSTRAINTS AND OPPORTUNITIES

A number of technological trends are threatening to impede innovation in computing systems, while others are creating new innovation opportunities. We start with three technology constraints that are driving up the cost of growth in computing systems.

## 3.1. CONSTRAINTS

A decade into the 21st century, computing systems are facing a once-in-a-lifetime technical challenge: the relentless increases in raw processor speed and decreases in energy consumption of the past 50 years have come to an end. As a result, the whole computing domain is being forced to switch from a focus on performance-centric serial computation to *energy-efficient parallel computation*. This switch is driven by the higher energy-efficiency of using many slower parallel processors instead of a single high-speed one. However, existing software is not written to take advantage of parallel processors. To benefit from new processor developments, software developers must re-design and rewrite large parts of their applications at astronomical cost.

### EXAMPLE OF RETHINKING AN APPROACH: EASY TASK PARALLELISM WITH “UNCERTAINTY REDUCTION”

Previously, for increasing the accuracy of large modeling code – like weather forecasting – the trend was to use a finer grid. But another solution is not to reduce the size of the calculation mesh, but to repeat the computation with a slightly different set of input data in order to reduce the uncertainty. This replication of code with other input parameters is easier to

implement because each task is independent of the other (no strong dependencies between different instances of tasks on the machine) and works well on data servers (as long as each node has enough memory to store the code and its parameters).

(From <http://www.semico.com/press/press.asp?id=299>)

Yet even the shift to universal parallelism is not enough. The increasing number of components on a chip, combined with decreasing energy scaling, is leading to the phenomenon of “dark silicon”, whereby the chip’s power density is too high to use all components at once. This puts an even greater emphasis on efficiency, and is driving chips to use multiple different components, each carefully optimized to efficiently execute a particular type of task. This era of heterogeneous parallel computing presents an even greater challenge for software developers. Now they must not only develop parallel applications, but they are responsible for deciding what types of processors to use for which calculations.

Tackling these issues requires addressing both hardware and software challenges. We must design energy-efficient systems with the right mix of heterogeneous parallel components and provide developers with the tools to effectively leverage them. Without either development, we will be unable to continue the computing systems growth that has dramatically changed our society over the past 50 years. Accomplishing this will require a global reassessment of how hardware and software interact.

### 3.1.1. HIGH PERFORMANCE HARDWARE BLOCKED BY THE FOUNDRY COSTS

The cost of the advanced process technologies required to continue Moore’s law is growing exponentially. This cost is doubling every four years, making the continuation of Moore’s law not only a technical challenge but also an economic and political one. The chart below shows the evolution of foundries for the past technology nodes. For every new technology node, there are fewer foundries left, and for the most advanced nodes, none of them are fully European.

On top of that, the non-recurring engineering costs of designing new chips are also growing rapidly due to the complexities of increasing levels of functional integration and the difficulties of managing smaller transistors. From these trends, it is clear that investing in a new fab is only worthwhile if it can produce extremely large volumes, i.e., billions of chips. The growing non-recurring engineering costs of chips will similarly lead to the need to sell larger volumes to break even. These economic realities will lead to a decreasing diversity of chips and a dramatically contracting selection of fabrication vendors.

	130nm	90nm	65nm	45/40nm	32/28nm	22/20nm
Altis Semiconductor	●					
Dongbu HiTek	●	●				
Freescale	●	●				
Fijitsu	●	●	●	●		
GlobalFoundries	●	●	●	●	●	●
Grace Semiconductor	●	●				
IBM	●	●	●	●		
Infineon	●	●	●			
Intel	●	●	●	●	●	●
Renesas (NEC)	●	●	●	●		
Samsung	●	●	●	●	●	●
Seiko Epson	●	●				
SMIC	●	●	●	●		
Sony	●	●	●			
STMicroelectronics	●	●	●	●	●	
Texas Instruments	●	●	●			
Toshiba	●	●	●	●		
TSMC	●	●	●	●	●	●
UMC	●	●	●	●	●	

Production capabilities in logic CMOS technology for main semiconductor manufacturers (2011) from HIS iSuppli  
Source: IHS iSuppli

### 3.1.2. THE POWER CONSUMPTION COST

A major paradigm shift is now taking place. While Moore’s law will keep its pace, continuing to double the transistor density, it will only allow for a minor increase of clock frequency and decrease of power dissipation per transistor. This evolution signifies the end of “Dennard scaling” and the resulting impact of “Dark silicon”.

#### • End of Dennard scaling:

Dennard scaling implies that voltage, current, and capacitance scale proportionally to the linear dimensions of a transistor, as does maximum frequency. Under this scaling, power consumption will be proportional to the area of a transistor and power density per device is nearly constant. This property implies that as we shrink transistors they will consume less power, thereby allowing the total chip power to remain constant even as we increase the

## MOORE'S LAW" STILL ON FOR THE NEXT DECADE

With the new silicon technologies like FinFETs and FDSOI, we can produce bulk transistors beyond the 20nm node. FinFETs are a 3D technology for controlling the conduction channel of the transistors. FDSOI is a 2D technology that allows transistors to run at frequencies up to 30% faster than bulk CMOS, while being 50% more power efficient with lower leakage and a much wider range of operation points down to lower voltages.

For example, the same ARM Cortex A9 processor can run at 150 MHz at 0.6V in 28nm Low Power bulk technology and up to 800 MHz in 28nm FDSOI with body biasing. The current semiconductor roadmap promises a 40% increase in performance from 28nm FDSOI to 14nm FDSOI (2014) and more than a 30% performance increase from 14nm to 10nm FDSOI (2016).

*(data from ST Microelectronics and CEA-Leti)*

FinFET technology also offers speed and energy improvements: for example, an ARM core design running at 3GHz consumes 3W

in 28nm bulk (28HPM), 2W in 20nm (2oSOC) and only 1.2W in 16nm FinFET. The 16nm FinFET can operate from 0.5V to 0.8V, with only a 50% speed reduction at 0.5V, compared to a 75% reduction for 20nm bulk. The tentative FinFET roadmap from TSMC will have 16nm devices by the end of 2013, 10nm by 2016 and 7nm by 2019.

*(data from TSMC)*

To produce these smaller features, the fab industry should move to EUV (Extreme ultraviolet lithography). This represents a drastic change in the micro lithographic technology. EUV must be done in a vacuum, and all optical elements, including the photo mask, must make use of defect-free multilayers that act to reflect light by means of interlayer interference instead of transmission of light.

To move to even more advanced nodes, e.g. 5nm, new techniques will be needed to control the transistor gate on all four sides, for example by using nanowires.

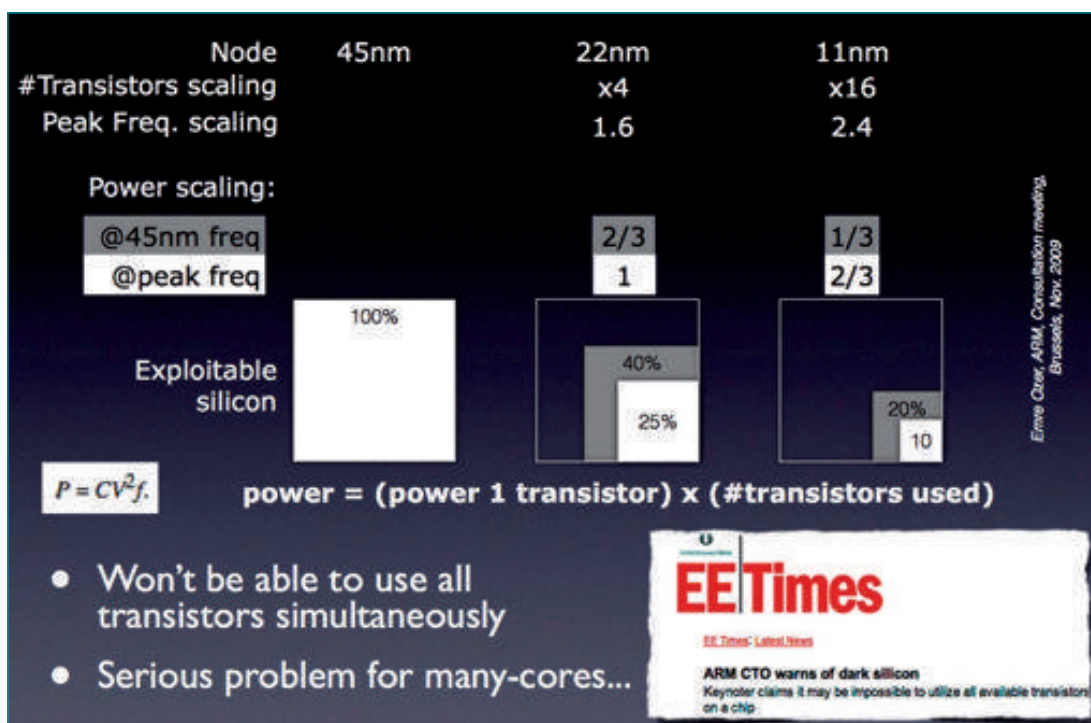
Unless some unexpected problems arise, it is generally expected that Moore's law will continue for the next 10 years.

number of transistors. This rule held for many decades, but ceased in the early 2000s: voltage scaling and frequency increases dropped significantly, and the power density started growing with the linear dimension of transistors, therefore limiting the number of devices that can dissipate power on a chip.

### • Dark Silicon

Even if it will still be feasible to pack more devices on a chip, the

power dissipation of each device will not be reduced accordingly due to the end of Dennard scaling. Since we are already pushing the limits of power dissipation/consumption on a chip, it will no longer be possible to power on all transistors on future chips simultaneously. Doing so will either dissipate too much heat or consume too much energy. This inability to turn on all of the transistors on a chip is known as "Dark Silicon".



Dark Silicon (From ARM) [Aitken1], [Esma1]

• **Communication**

Communication, or moving data, is a key problem in terms of system efficiency and power. The energy required to do a computation in modern chips is dwarfed by the energy required to move the data in and out of the processor. This is forcing developers, for both hardware and software, to figure out how to change their systems and software to keep data as physically close as possible to the computing elements. Similar issues arise at larger scales in high performance and data center computing, where the energy required to move data between thousands of

machines is often much higher than that required to compute the desired answer.

In the picture below, we can see that the energy of a 64-bit datapath has the same order of magnitude as moving data from one side of the chip to another, but sending or receiving data from an external interface is several orders of magnitude higher in terms of its energy requirement. Coping with the energy requirements for moving data requires developments at both the hardware and software level.

**ENERGY FOR COMPUTING VERSUS ENERGY FOR MOVING DATA**

In 22nm, swapping 1 bit in a transistor has an energy cost of roughly:

$$\sim 1 \text{ attojoule } (10^{-18} \text{ J})$$

Moving 1 bit on the silicon costs approximately:

$$\sim 1 \text{ picojoule/mm } (10^{-12} \text{ j/mm})$$

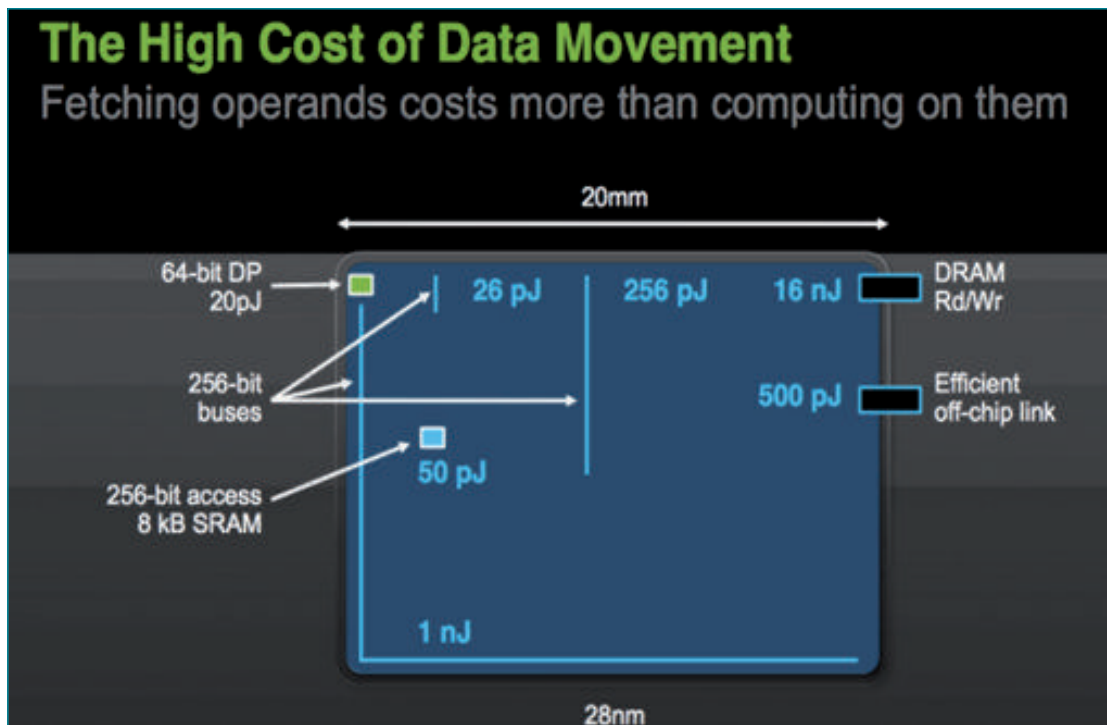
Moving data  $10^9$  times per second (1 GHz) on silicon has a cost of:

$$1 \text{ pJ/mm} \times 10^9 \text{ s}^{-1} = \sim 1 \text{ milliwatt/mm}$$

$$64 \text{ bit bus @ 1 GHz: } \sim 64 \text{ milliwatts/mm (with 100\% activity)}$$

$$\text{For 1 cm of 64 bit bus @ 1 GHz : } 0,64 \text{ W/cm}$$

On modern chips, there are kilometers of wires on chip, and even with a low toggle rate, this leads to several Watt/cm<sup>2</sup>



Source: Bill Dally, *To ExaScale and Beyond*. [www.nvidia.com/content/PDF/sc\\_2010/theater/Dally\\_SC10.pdf](http://www.nvidia.com/content/PDF/sc_2010/theater/Dally_SC10.pdf)

## THE CHALLENGES OF DATA-CENTRIC COMPUTING

- 1) Moving data becomes the major problem in architecture in general, and in particular for data servers.
- 2) The second limitation of large data centers is the energy consumption, because of the cost of the supply of electricity. Any growth in computing power is only possible at the same energy consumption (“*iso-energy*”): typically the considered budget is about. The “power wall” for a data center comes from the physical difficulty of providing more power to the facility in the form of power lines and transformers.
- 3) The increase of computational performance is a third priority after the power issues and data access.

For data servers, one current “hot” challenge is the one of “big data” and especially of “dark data” (data stored in data servers that are never replayed). Recent estimates assess the ‘dark data’ at 95% of the stored data. They clutter the system resources and limit the time for accessing data. The problem is often exacerbated by the distribution of the data that prevents the optimal use of communication links.

Data servers in the future will no longer be “Computation Centric”, but rather “Data Centric” aka “*Data Centric Data Center*”: it will be easier to move a computation to the data than to move the data to a computer. The system should bring the algorithm, the computation close to the data.

To address this we must rethink existing architectures and follow a holistic interdisciplinary approach to profoundly change the ecosystem.

### 3.1.3. COMPLEXITY: THE ACHILLES HEEL OF SOFTWARE

Multi-core and heterogeneous systems have been imposed upon software developers due to technological constraints. And while these systems have phenomenal “peak” performance, it is incredibly difficult to get real-life applications to approach that performance.

It will ultimately be necessary to make parallel programming as simple as sequential programming if we are to take advantage of these processors across all domains. However, most contemporary parallel languages are primitive and low-level: they are for the parallel computers what machine language was for the first sequential systems. They are cumbersome and require intricate knowledge of the execution model of the machine to fully optimize the performance. Most critically, they do not provide portable performance between different hardware platforms. Although there has been progress over the last decade, the challenge of enabling cost-effective, efficient multi-core

programming is still a primary concern in computing systems. Adding heterogeneity, dark silicon, and scaling to high core counts only makes the challenge more daunting.

Another challenge is at a much higher level of abstraction: programming a network or system of computing systems (as it is done in data centers, and between mobile devices and sensing devices). There are already some successful models like map-reduce, but these are only applicable to particular types of workloads (the so-called scale-out workloads). For less regular algorithms, there is no ready-made solution to map them onto a complete parallel system. The current best practice is to do so manually, but this requires extensive expertise on the behavior of the different computing platforms and of their interaction. Debugging such systems is a nightmare due to a dearth of advanced debugging tools for heterogeneous, parallel and distributed systems.

## 3.2. OPPORTUNITIES

### 3.2.1. ENTERING THE THIRD DIMENSION

A logical continuation of Moore’s law is to stack dies and interconnect them in the vertical direction. Several techniques to enable high-bandwidth interconnect between dies already exist, such as Cu-Cu interconnect, Through Silicon Vias and optical connections. Stacking dies creates several opportunities:

- It enables us to change architectures by physically reducing the distance between different devices by stacking. For example, we can place memories and processors (for many-core architectures) or sensors and processing (e.g. intelligent retinas) on top of each other to increase integration of different devices.
- It allows us to combine different technologies in one package, meaning that not all the dies in a package need to be produced in the same technology node. This can extend the lifetime of existing fabs by, for example, producing a sensor chip in an older technology and mounting it on top of a modern processor die.
- Through different combinations of dies, we may be able to regain the chip diversity lost due to the increasing costs of chip design, and particularly to the cost of semiconductor chip fabrication plants that doubles every four years (Rock’s law). By reusing dies in different packages, the die volume of one design will increase, thereby lowering the cost, while simultaneously providing for more differentiation through different stacking combinations.

3D Stacking also has its own set of challenges that must be overcome to make it economically viable at a large scale. These hurdles include reliably connecting the dies, dissipating the heat from the different dies, and distributing the functionality among the different dies. To increase yield, technologies such as die-to-

wafer bonding or die-to-die bonding have to be improved compared to the current wafer-to-wafer approach, which is easier to achieve but requires the yield of each wafer to be very high.

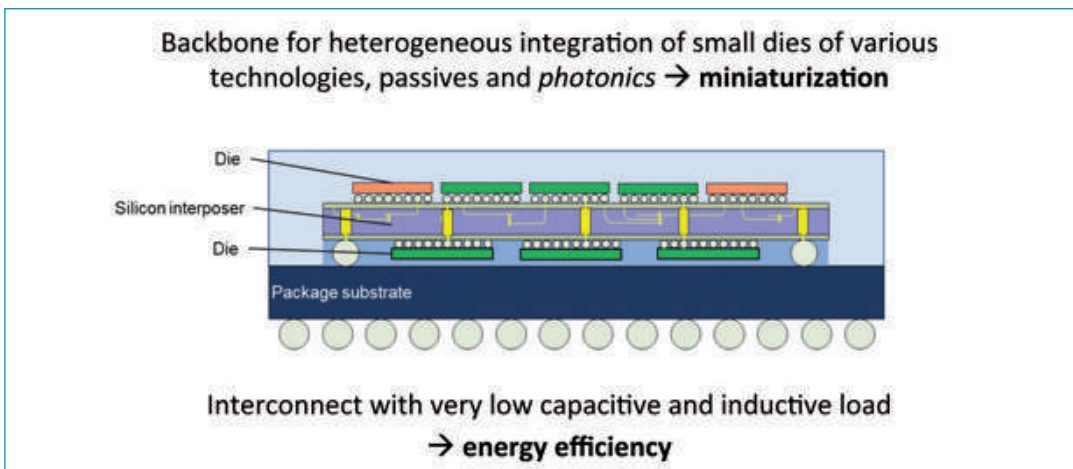
Before the technology matures for full 3D stacking, systems using silicon interposers (also called “2.5D” solutions) will become mainstream. These systems use a silicon substrate and metal interconnect to create a highly integrated circuit board that can connect multiple dies together. These dies can be directly connected to both sides of the interposer (using Cu-Cu interconnect for example, or using TSVs if necessary). The technology is less demanding than full 3D stacking while enabling high levels of integration at reasonable costs.

Silicon interposers can also be active (i.e. integrate active functions, such as Input/Outputs or converters). By integrating multiple dies, each die can be developed in an optimized technology. This

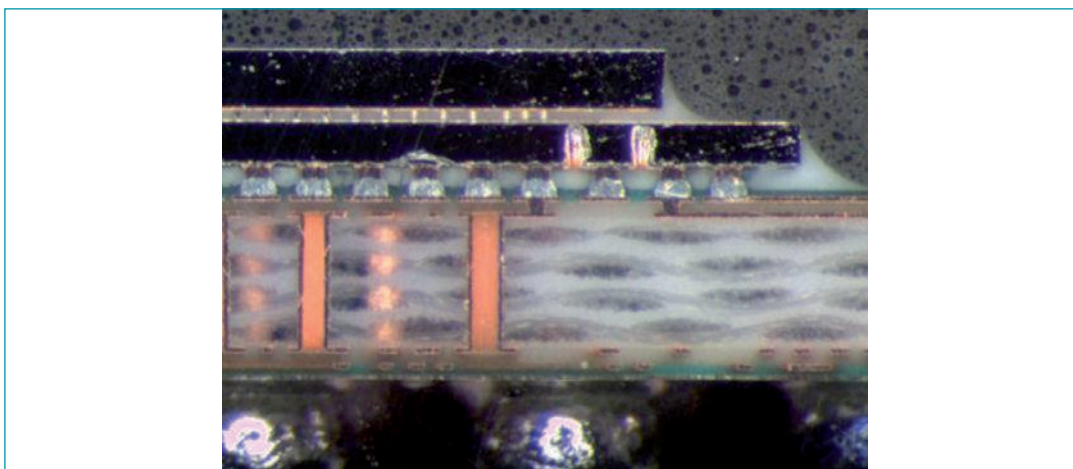
allows for cheaper dies (using only advanced processes for the most performance-critical dies) and more diversity (by integrating circuits that require optimized processes, such as analog, digital, photonics, and sensors). Silicon interposers do not require the most advanced technology nodes, and therefore can be built in existing fabs.

Entering the third dimension may bring solutions for the challenges of power (by reducing the length of interconnect), diversity (making a “System on Interposer” composed of functions where each of them is integrated into its optimal technology), and cost for building highly complex monolithic SoCs.

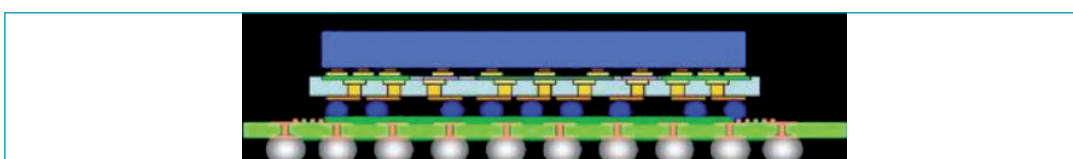
Silicon interposers are also promising for integrating silicon photonics, and therefore enabling optical interconnects between systems or dies.



2.5 D silicon interposers



3D stacking



Photos: STMicroelectronics & CEA-Leti: Multi-die stacking using Copper-pillars and TSVs

### 3.2.2. SILICON PHOTONICS

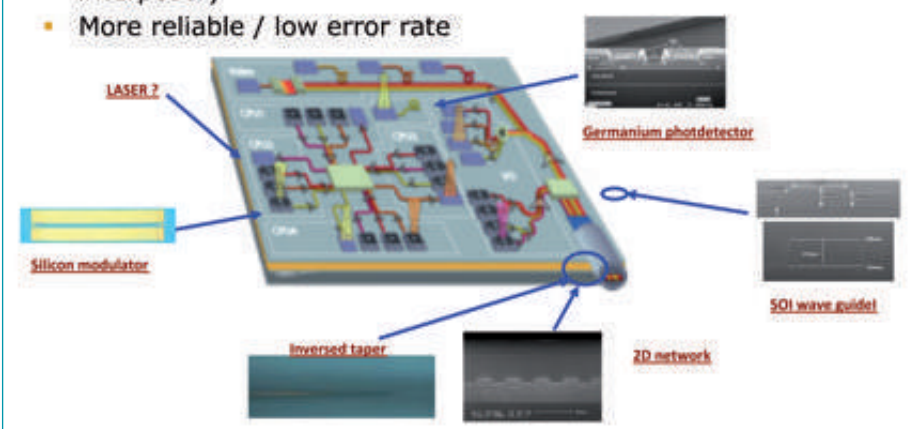
Another technology is silicon photonics, which holds the promise of lower communication energy cost, higher bandwidth and low manufacturing cost in existing fabs. The technology is compatible with existing silicon technology but even if it doesn't require the latest nanometer technologies because of the micrometer size of the various devices, the current precision for building the devices is in the nanometer range. Transforming the electrical signal into an optical one is also power consuming, mainly if serial conversion

is required (leading to very high frequencies), therefore optical systems offering multiple wavelengths might be preferable from an energy point of view. Current optical systems are more at the metascale level (optical fibers connecting Internet devices across countries or oceans, or connecting racks in data centers), but they are also promising for connecting boards or even chips together with a very high bandwidth.

Silicon photonics is a solution that can reduce the problem of interconnects, which require high bandwidth and power.

### CMOS photonic: integration of a photonic layer with electronic circuits

- Use of standard tools and foundry, wafer scale co-integration
  - -> Low manufacturing cost
- **Lower energy** ( $\sim 100$  fJ/bit), (wire:  $\sim 1$  pJ/mm) -> less heating
- High bandwidth (10 Gbps), Low latency ( $\sim 10$  ps/mm)
- High integration
- Can also be used for on-chip clock distribution (using e.g. passive SOI interposer)
- More reliable / low error rate



Summary of advantages of Silicon Photonics from Computing Systems

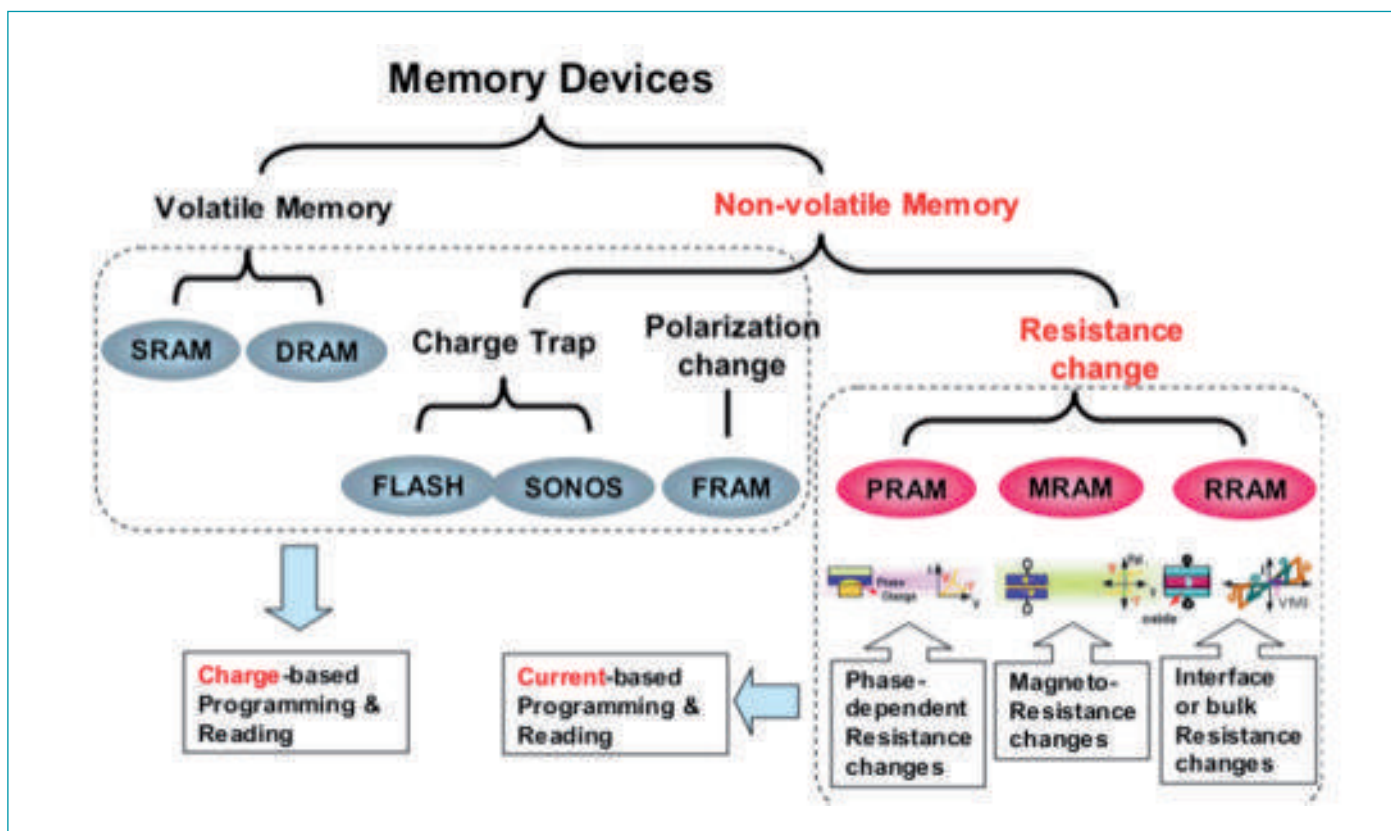
### 3.2.3. WIRELESS CONNECTIVITY

Wireless communication bandwidth is now approaching the speed of low-cost hard disks. This is a game changer as it means that – given affordable communication cost – it is no longer required to locally store data in smart devices, but instead the devices can directly connect to the internet and remote storage. The “cloud” is born from this evolution. This evolution creates opportunities of new system design, shared address spaces over wirelessly connected devices, etc. It can also avoid replication of data on multiple mobile devices (pictures archives, music libraries). Of course, the limitations (constant connectivity

required, power dissipation) need to be taken care of when designing the complete system.

### 3.2.4. EMERGING MEMORY TECHNOLOGIES

Among the broad scope of emerging technologies, non-volatile memories offer interesting properties that can be used in future computing systems. Those novel memory devices can be classified into various categories depending on the physical processes implementing the memorization process. See the following figure for a brief taxonomy.



A brief taxonomy of memory device technologies

Of all technologies depicted above, MRAM and PRAM devices are the most advanced from an industrial point of view. MRAM was first introduced in 2004 by Freescale, and with today's refinement in technology, it has allowed Everspin to introduce a 64 Mb ST-MRAM chip. PRAM is currently in active development and industrial products are on the way with Micron just announcing a 1Gb PRAM module for mobile applications. Resistive ram (RRAM) comes in many technological flavors such as Conductive-Bridge or Oxide Ram; they are actively developed by industrial companies such as Adesto technologies or Hynix and may prove to be the best option for future very dense arrays of non-volatile memory.

#### Emerging Non Volatile memories for digital design

In Reconfigurable logic circuits, non-volatile memories can be used to implement Look up Tables or to implement switch box matrices as demonstrated on a small scale by HP [Xiao9]. This later work was also one of the first demonstrations of the co-integration of titanium dioxide memristors on top of an industry CMOS process. More original approaches to implement sequential or stateful logic circuits using memristive devices have been proposed [Borgh10].

Since novel resistive devices are dipoles (i.e. they are two terminals devices), for which resistance can be controlled by applying voltage levels, the idea to organize them in crossbars is pretty natural with the expectation of getting ultra-dense arrays of

non-volatile memory bits. Although using the crossbar scheme with devices of nanometric sizes can provide dramatic memory densities, it is not without a lot of problems to solve: select logic, sneak paths, and process variability, to name a few. However, a convincing demonstration of a 2Mb resistive device crossbar without selection logic has been recently unveiled by Hynix and HP [Lee12]. This work shows that putting those devices to practical use in digital design would be possible at the price of rethinking the whole architecture of the memory plane. Up to now, practical demonstrations have targeted either reconfigurable logic or full arrays of memories but a very promising possibility for non-volatile memories could be to use them in the design of processor register files with the aim of building instant-on/instant-off CPUs.

For digital applications, the key selling points of those novel technologies is their potential to yield ultra-dense arrays of non-volatile devices with possible CMOS-compatible processes. Using non-volatile memories at key points in a computing architecture can dramatically reduce its overall power consumption.

#### Emerging technologies for Neuromorphic computing

The relation between the hysteretic voltage-current diagram of Memristive technologies and that of an artificial synapse has been put forward by several authors: Leon Chua in his seminal paper on the memristor [Chua71] and more recently the HP labs

## “ELECTRONS FOR COMPUTE, PHOTONS FOR COMMUNICATION, IONS FOR STORAGE AT NET ZERO ENERGY”

(from Partha Ranganathan [Partha12])

Besides the current technology for compute nodes, evolving with FinFETs or FDSOI, which rely on electrons, two other elements are key enablers for future compute systems: photonics for interconnects and non-volatile RAMs for storage.

Light must ensure a 30x increase in bandwidth at 10% of the power budget of copper.

The current problem of photonics, with “ring resonators”, is that energy is consumed even if the device is not used. There is an important fixed cost. Therefore the optical link must be shared (communications multiplexed on the same channel).

Bandwidth to memory is very important, and it is a good field for improvement, but there is the problem of inheritance for the

DRAM industry, with a low possibility of impacting it with new architectures, except when it will hit a wall. This will then be a good time for emerging technologies such as the NVRAMS and memristors. Non-volatile memory concepts, with their properties of addressing at the byte level, will induce new approaches to define efficient file systems that are not based on the concepts of blocks of data. This byte access will also change the structure of the data cache.

However, the memristors will not resolve the access latency problem (but perhaps will decrease the requirement for L2 caches). On the other hand, for reasons of low dissipation, compute nodes will work near the threshold voltage with working frequencies in the sub-GHz range. Simultaneous multithreading can be used to hide memory latency. Memristors can be used to save the state of the processor when task switching. It would take about 100 threads to hide the latency of memory accesses, which is accessible for data server applications.

team [Strukovo8]. Furthermore, it has been proposed that by combining the symmetric thresholds of a memristive device together with a spike based coding, the induced resistance change could result in the implementation of a learning scheme known as STDP (Spike Timing Dependent Plasticity) [Snidero8] potential scalability, and inherent defect-, fault-, and failure-tolerance. We show how to implement timing-based learning laws, such as spike-timing-dependent plasticity. An important body of work has been triggered in an effort to implement dense arrays of synapse-like devices using emerging memory technologies.

It is interesting to note that technologies such as PRAM, MRAM or RRAM offer many opportunities both in the exploration of new digital processor architectures and in the investigation of disruptive computing paradigms and coding (spikes) as exemplified by neuromorphic computing.

For the reader interested in this burgeoning topic, a valuable review on emerging memory technologies can be found in [Pershin1].

### 3.2.5. STOCHASTIC/APPROXIMATE COMPUTING

As a lot of data that will be processed comes from the natural world, exact processing with a lot of accuracy is not always necessary. Several approaches and techniques can be used to reduce the power, or increase the density, of processing engines, taking advantage of “less accurate” processing.

#### Probabilistic CMOS:

This approach, pioneered by K. Palem [Cheemo5], consists of using standard CMOS, but lowering the voltage of the transistors. As a consequence, the energy is significantly reduced but the transistor provides the correct output only with a certain probability. However, large classes of important algorithms (e.g., optimization algorithms) are compatible with such a probabilistic computing medium. First prototypes show very promising energy gains.

#### Approximate computing:

The power consumption of CPUs and memory systems has traditionally been constrained by the need for strict correctness guarantees. However, recent work has shown that many modern applications do not require perfect correctness [Samp1]. An image renderer, for example, can tolerate occasional pixel errors without compromising overall quality of service. Selectively-reliable hardware can save large amounts of energy with only slight sacrifices to quality of service (from [Sampa]).

#### Graphene or Graphyne devices [Schirber]:

Other potentially disruptive technologies are emerging, such as graphene (Physics Nobel Prize in 2010) transistors, which seem capable of increasing clock frequency beyond the capabilities of silicon transistors (in the 100 GHz to THz range). Currently such transistors are significantly bigger than silicon transistors, and only limited circuits have been implemented. Their application scope is mainly fast and low-power analog signal processing, but research on graphene transistors is still in its infancy, and this technology should be carefully monitored by the computing

industry. Computing might shift back to the frequency race instead of parallelism if complex devices running with low power at 100 GHz are possible [Manchester].

### 3.2.6. NOVEL ARCHITECTURES

Previous approaches use an explicit declaration of how to perform tasks (typical *imperative programming*), but we can think of paradigms where, instead of instructing the machine on how to perform its tasks, we only specify the goal(s) or the objectives. *Declarative programming* (like database query languages - e.g., SQL, regular expressions, logic programming, and functional programming) falls in this category, together with other approaches like using Neural Networks. Those approaches are promising to cope with the complexity of programming large-scale parallel and/or distributed systems. Most of them can be easily mapped to parallel systems.

Gaining inspiration from the brain is one way to improve our computing devices and to progress beyond Moore's law. As high-performance applications are increasingly about intelligent processing, the application scope of neural networks becomes very significant. More importantly, as the number of faulty

components increases, hardware neural networks provide accelerators that are intrinsically capable of tolerating defects without identifying or disabling faulty components, but simply by retraining the network. They will not replace a complete system, but used as accelerators they could decrease the computational load of classical systems in an energy efficient approach, mainly if implemented with new devices.

Since their introduction memristors have been recognized for their applicability as synapse-like elements and the possibility for a direct implementation of the STDP (Spike-Timing-Dependent Plasticity) learning rule [Snider08]. Similarly, organic devices have been demonstrated to have synaptic properties [Alibart10] and recently phase change memory was proposed for the same purpose [Kuzum11]. Beyond the fact that the STDP learning strategy implies a paradigm shift in information coding (from state based to event based), it promises easier implementation of very dense arrays and more energy efficient systems for qualified classes of applications (e.g. sensory data processing, vision). Big research initiatives in the field of neuromorphic computing are currently under way, such as the DARPA funded SYNAPSE project coordinated by IBM [SyNAPSE] and in Europe several ICT projects with FP6-FACETS/FP7, Brainscale being the most famous.



# THE POSITION OF EUROPE

In this section, we give a SWOT-analysis for computing systems in Europe. It is an important input for the recommendations, as we want to build on strengths to remediate the weaknesses.

## 4.1. STRENGTHS

### 4.1.1. STRONG EMBEDDED ECOSYSTEM

The European computing industry has a strong embedded ecosystem spanning the entire spectrum from low power VLSI technologies to consumer products. Companies such as ARM, Imagination Technologies, STMicroelectronics and ST-Ericsson are leaders in providing semiconductor processing elements and IP for embedded systems. Imec and CEA-Leti are leading the development of semiconductor technology, while ASML is the leader in photolithography equipment. Large end-user European companies have a strong market presence internationally in areas such as automotive (Volkswagen, Renault-Nissan, Peugeot-Citroën, Fiat, Daimler, BMW, Volvo), aerospace and defense (Airbus, Dassault, Thales, Saab, Barco), telecommunications infrastructure (Nokia, Ericsson, Intracom), telecommunications operators (Deutsche Telekom, Telefónica, Orange), system integrators (Thales, Siemens), and software services (SAP). These companies are all globally competitive and work on the forefront of system design and implementation.

These larger players also rely on a thriving community of SMEs that strengthen the technical and innovative offers in the market. This strong embedded ecosystem creates a fertile environment for innovation, better integration and worldwide leadership. The SME-ecosystem also spans the complete spectrum.

There are many fabless semiconductor companies like Kalray, Recore Systems, Think Silicon, Clearspeed Technology, and NovoCore Ltd. These are supplemented by companies that build tools and carry out consulting for these tools like CAPS Enterprise, ACE, Maxeler Technologies, Modae Technologies, Vector Fabrics, Codeplay, CriticalBlue, Ylichron, and Leaff Engineering. Some companies are active in the real-time and safety critical domain like OpenSynergy, Sysgo, Rapita Systems and Yogitech; others do formal verification like Monoidics, or security like Herta or INVIA. We have also seen a decline of the previously big and vertically integrated European companies like Philips, Siemens, Thomson. Unlike companies like Samsung, they “deverticalized” by spinning out their different business, for example their semiconductor businesses. The resulting companies, NXP and Infineon, reduced the scope of their activities and their ranking dropped amongst major semiconductor players.

### 4.1.2. PUBLIC FUNDING FOR R&D AND TECHNOLOGY TRANSFER

Europe benefits from a large pan-European centralized research program through the Framework Programmes for Research and Technological Development. If used properly, these serve as a powerful tool to direct community-wide research agendas and a strong incentive for diverse groups of companies and research institutions to work together across political and cultural divides.

The value of the resulting consortia generally does not lie in bringing a particular new product or technology to market in the context of a project, but rather in investigating the potential of a new technology in terms of applications, business opportunities for individual partners, and creating new ecosystems and markets.

### Exploratory work, basic research

Basic research and exploratory work are high-risk activities for industry, and also harder to fund for universities than more applied research. At the same time, such fundamental innovation is one of the activities with the highest possible return on investment due to its potential to create whole new markets and product niches. It is therefore of paramount importance to strongly support such undertakings, and public funding is an excellent way to lower the financial risk by steering them towards selected domains and research directions.

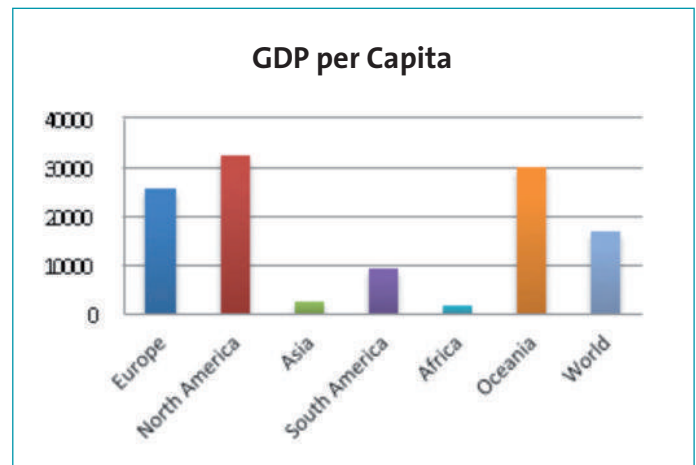
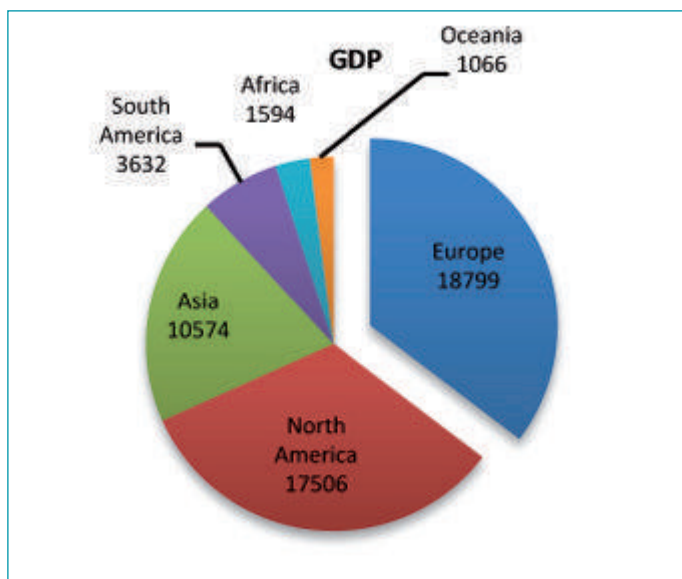
### Bootstrapping new business ecosystems

A product by itself seldom creates a new, large market. The real value of many products lies in the fact that they can be used in many different ways, combined with many different other products, integrated into existing workflows. These so-called network effects, whereby value increases through third party products and services making use of and interoperating with the original product, are particularly important in today's connected world.

Enabling such network effects requires setting up ecosystems, which often require expertise from many different horizontally specialized players: hardware design, development tools, programming models, drivers for third party hardware, operating system support, etc. Publicly funded projects are excellent drivers to bring all these players together and help fund the necessary support and integration. The resulting, often at least partly open, environments encourage the creation of startups by offering a stable foundation to build on.

### 4.1.3. ONE OF THE BIGGEST MARKETS

According to the World Bank, the world GDP is distributed over the continents as follows:



Continent	GDP	GDP/Capita
Europe	18799	25467
North America	17506	32296
Asia	10574	2539
South America	3632	9254
Africa	1594	1560
Oceania	1066	29909
World	53175	16837

From this table is clear that Europe is the biggest economy of the world with a lot of potential for growth in recently joined member states, and large neighboring markets (Russian states, Middle East, Africa). With a high GDP per Capita, there is a lot of money to spend on innovative products and services.

### 4.1.4. GOOD EDUCATION

From an educational perspective, more than 200 European universities are rated among the top 500 universities in the Shanghai Jiao Tong University ranking [Arwu2012]. This is more than in the United States of America (154 universities). Europe benefits from a very strong educational environment and a highly competitive undergraduate and graduate educational system. The ongoing bachelor-master transformation will further strengthen the European educational system.

## 4.2. WEAKNESSES

### 4.2.1. EUROPE IS FULL OF HORIZONTAL SPECIALIZATION

Very few European companies offer a completely integrated vertical stack that they control from top to bottom (services, software, and hardware). This sometimes makes it hard to compete with vertically-integrated US and Asian giants that can easily lock-in customers from the point that they buy a particular device, or at least entice customers into using the manufacturers' own software and services.

#### 4.2.2. LOSS OF COMPETIVENESS IN SOME DOMAINS

European players own fewer patents on emerging standards. And while patents on standards are a controversial topic, a constant factor is that whoever owns patents that are part of widely accepted standards can make a lot of money, and others have to pay. Looking at previous MPEG video and audio compression standards and comparing them to upcoming ones such as H.265, there is a clear shift in patent ownership towards Asian companies.

Another weakness is the decreasing number of state-of-the-art foundries in Europe, which could result in a decrease of know-how about making advanced complex circuits and a shift of fabrication outside Europe. In the future, Europe might depend on strategic decisions made outside of Europe.

#### 4.2.3. BORDERS AND DIFFERENT LANGUAGES

Language and cultural diversity in Europe are handicaps to attracting bright international students to graduate programs outside of their home country. Lack of command of English by graduates in some countries also greatly hampers international networking, collaboration, and publication. The fact that multiple scripts are being used throughout Europe (Latin, Greek, and Cyrillic) makes integration also more difficult.

#### 4.2.4. LACK OF VENTURE CAPITALISTS

The lack of a venture capitalist culture contributes to the brain drain. It is much harder for a university or PhD graduate to start a company in Europe than in the United States. Yet even with venture capital, the attitude of personal investment and identification of startup employees with the fate of their company that drives Silicon Valley is largely absent from the European work ethos and organized labor agreements. On top of this, bureaucracy and administrative procedures in some countries are preventing or killing several new initiatives.

#### 4.2.5. WEAK ACADEMIA-INDUSTRY LINK

European computing research is characterized by a weak link between academia and industry, especially at the graduate level. Companies in the United States value PhD degrees much more than European companies, which often favor newly graduated engineers over PhD graduates. This leads to a brain drain of excellent computing systems researchers and PhDs trained in Europe to other countries where their skills are more valued. As a consequence, some of the successful research conducted in Europe ends up in non-EU products or does not make it into a product at all.

### 4.3. OPPORTUNITIES

#### 4.3.1. COST EFFECTIVE CUSTOMIZATION

Building a state-of-the-art fab is very expensive and it requires a 95-100% load to make it profitable. In previous years, most of the European Silicon makers have become fab-less or “fab light” (NXP, Infineon, ..), and don't plan to keep up with the most advanced technology processes. Building modern monolithic SoCs requires a \$0.5-\$1B investment or more, which can only be amortized with millions of pieces per year. This is very difficult for Europe, but we have seen in section 3.2.1 that 2.5D technology (Silicon Interposers) could leverage current European fabs while offering a way for European companies to differentiate and integrate. It can also leverage the know-how of Europe in embedded systems, MEMS and other integrated sensors and devices for building advanced composite circuits with an affordable design-start cost (\$10M vs. \$1B). Some roadblocks are still present: for example, building interoperable European components for assembly on a silicon interposer, etc.

#### 4.3.2. LEVERAGING FREE/CHEAP/OPEN INFRASTRUCTURE

Open Source and Free Software has spread to virtually all corners of the computing market. This evolution can be, and in fact already is, leveraged by companies in order to distribute the cost of implementing and maintaining basic infrastructure that is not part of their core business. They may need the resulting tools themselves, or may want to build ecosystems around their products. Such ecosystems themselves offer many opportunities for starting new businesses, both in terms of directly using the software (hosting providers, social media, wireless routers, etc.) and in terms of offering consulting services to adapt such software for particular purposes.

This evolution is also starting to take place in areas like content and services (e.g. the UK-based OpenStreetmap) and hardware (e.g. the UK-developed Raspberry Pi).

#### 4.3.3. SOCIETAL CHALLENGES

Paradoxical as it may seem, several challenges that society is facing are also huge opportunities for research and industry in computing systems. For example, the aging European population will require the development of integrated health management and support systems that enable people to live at home for a longer time. Europe's national health systems are far better placed to take advantage of and coordinate such efforts than those of the United States or China. European expertise in low-power and embedded systems and its SME ecosystem is an asset for tackling other grand challenges such as the environment, energy and mobility. Experience in mission critical systems gives

Europe a competitive advantage with the safety and security challenges that lie ahead in larger-scale consumer systems.

#### 4.3.4. CONVERGENCE

Disruptive technologies such as cloud computing and the convergence of HPC and embedded computing represent opportunities for Europe too. The trend towards more distributed environmentally-integrated cyber-physical systems could be beneficial to the European semiconductor industry, which has significant expertise in the wide range of required technologies.

#### 4.3.5. MICRO- AND NANO-ELECTRONICS

The recent move to classify micro- and nano-electronics as key enabling technologies [KET] for Europe creates significant opportunities for the computing systems industry in Europe. This move comes in response to Europe's semiconductor market share decreasing from 21% to 16% since 2000 [KET]. Such large and centralized research and development planning presents a great opportunity for computing systems research in Europe, bringing it on par with other technologies such as energy, aerospace and automotive technology. The resources available for developing pan-European research capabilities could be used to address several of the weaknesses mentioned above, in particular the lack of tools and hardware development expertise. The emphasis on academic collaboration with industry, in particular SMEs, can provide mentoring and bridges for researchers wishing to industrialize their results, if used properly.

### 4.4. THREATS

#### 4.4.1. COMPETING WITH FREE/CHEAP/OPEN INFRASTRUCTURE

##### Software: freely available software is “good enough”

When software that is available for free has functionality similar to commercial offerings, the latter can become a hard sell. Even if a commercially-supported version has more functionality, provides a better interface or comes with guaranteed support, it can be hard to convince current and prospective customers that these additional offerings are worth the financial premium at purchasing time. Moreover, in case the alternative is Open Source or Free Software, the ability to continue development should the original developers disappear is a hard to beat advantage.

##### Hard to sell development tools

Since development tools are a fundamental infrastructure required by many companies, much of the commercially sponsored Open Source and Free Software development happens in this area. Therefore the “freely available software” evolution particularly challenges companies whose core business consists of selling development tools. Moving towards adding value on top of free tools, and/or convincing customers that the added value of their own tool chain is worth the licensing cost, is far from trivial.

##### Services: Google is “good enough”

The wide variety of services offered for free by Google and others can cannibalize business and business opportunities for other companies. At the same time, companies that rely on these services to run their business are at risk due to the fact that such services often can be terminated at any time without notice.

##### Hardware: subsidized abroad

Developing hardware is considered strategic for some countries (e.g. China). Therefore some governments are directly or indirectly supporting development of microprocessors and other advanced digital circuits in order to gain strategic independence from external providers.

##### Open Source/Free Software licenses and business models

Not all Open Source and Free Software licenses are compatible with all business models. In particular the GNU General Public License (by design) can make it hard for proprietary software vendors to leverage code. Releasing the results from publicly funded projects at least under an additional license that makes such integration easier, such as the Modified BSD or MIT license, can help alleviate these problems. A parameterized solution may be useful, like the DESCAs templates for FP7 consortium agreements.

#### 4.4.2. FINANCIAL CRISIS

The financial crisis has caused many companies to focus more on their core business and reduce R&D spending. Several governments are also slashing education and public research funds in their efforts to balance their budget.



# CONCLUSION

Innovation in computing systems is going through turbulent times – the free lunch of automatic exponential performance increase is over. For almost a decade, the still-growing transistor budgets have been used to add more cores, but the tools have not been able to catch up with this evolution and common application development today is in essence still the same as a decade ago, while the complexity of modern systems has increased dramatically. Maintaining rapid growth in computing performance is key for tackling the societal challenges shaping Europe and assuring our global competitiveness in the future.

The traditional computing systems market is being replaced by a new market of smart *embedded* systems, *mobile* devices, and large-scale *data centers*, all converging to support *global-scale applications* that gather data from embedded systems and users, process it in large data centers, and control our environment or provide customized, timely information to millions of users through their mobile. The most successful global-scale applications are backed by global vertically-integrated companies offering a complete ecosystem.

Energy has become the primary limiting factor in the development of all systems, whether due to energy cost and cooling in large systems or due to battery life in mobile devices. This has led to the rise of parallel and heterogeneous devices that trade off increased *complexity* and incompatibility with existing software for higher efficiency, and to the appearance of “dark silicon”, whereby portions of a device must be shut off to stay within the power limit. The necessity to develop energy-aware devices and the ability to automate the optimization of applications for power efficiency has become a necessity across all computing systems.

Europe provides a strong embedded and low-power processor ecosystem. This includes many companies in hardware and software development for both the industrial and commercial sectors. However, Europe also suffers from a high degree of horizontal specialization, which makes it difficult for companies to amortize the costs of development across the product chain.

This roadmap updates the 2011 roadmap with new trends in market verticalization, global-scale computing, and the impact of hardware design costs, while continuing to emphasize the difficulties of achieving energy-efficiency and programmability across devices from the cloud to the mobile and embedded domains. From this analysis we identify *three strategic areas: embedded, mobile, and data center* and *three cross-cutting challenges: energy efficiency, system complexity, and dependability*.

## STRATEGIC AREAS FOR HORIZON 2020

### • Embedded systems

The traditional notion of an embedded system as a single-purpose device is rapidly changing as increased computing performance, connectivity, and closer interactions with the world bring additional functionality and demands. To take advantage of this potential we need to rethink system architectures and programming models; to optimize for energy, time constraints and safety; and to develop techniques to support portability of critical and mixed critical systems. Without such portability, the cost of certification for new computing platforms will prevent their uptake and limit our ability to leverage further advances.

### • Mobile systems

The shift from desktop PCs to mobile devices provides an incredible opportunity to rethink the human-computer interface and how we interact with technology. Innovations that provide more natural and immersive experiences will dramatically improve the utility and productivity of mobile devices. Such developments require collaboration across all levels of the computing system: from human interaction to image processing and data mining, down to system architecture for efficiency and performance. Further, as mobile devices become increasingly integrated in our public and private lives, we will need stronger guarantees of privacy and security.

### • Data center computing

As applications become global-scale, Europe has an opportunity to lead the global market for data center technology. To be competitive we must develop the capabilities to process “big data” without increasing cost or energy. These challenges include architectures for handling massive unstructured data sets, low-power server modules and standards, network and storage systems, scalable software architectures, and micro servers. At the same time we must integrate these developments with techniques to ensure security, privacy, and compliance, while providing large-scale reliability, availability, and serviceability.

## CROSS-CUTTING CHALLENGES FOR HORIZON 2020

### • Energy efficiency

Systems today are limited in their performance by power used or dissipated. This has led to the combination of many specialized (heterogeneous) processors to increase efficiency. Unfortunately, this has also increased the complexity of programming to the point where it is prohibitively expensive for many applications. On top of this, the energy cost of moving data now exceeds that of computing results. To enable power efficient systems we must

address the challenges of programming parallel heterogeneous processors and optimizing data movement, both for legacy applications and new computing modalities. We must also take advantage of the energy saving potential of new technologies, such as non-volatile memories, 2.5D and 3D integration techniques, new silicon technologies such as FinFETs and FDSOI, and new computing modalities such as stochastic and approximate systems and algorithms.

### • System complexity

Modern computing systems have grown to the scale where developers need to coordinate thousands of processors at once to accomplish complex tasks for large numbers of users and across massive data sets. To support this scale we need to develop tools and techniques to optimize for performance and ensure correct operation, while operating “at-scale”. On the hardware side, chips have become enormously more expensive to design, verify, and produce. Today’s cutting-edge technologies are so expensive that they are only affordable for devices that sell 10-100 million units. This cost limits product differentiation and makes market entry for new ideas extremely difficult. To overcome this we need to investigate new integration techniques that enable high levels of integration and differentiation without the cost of cutting-edge fabrication.

### • Dependability

Computing systems are involved in ever growing parts of our lives, from providing intelligent control for our transportation to keeping track of our friends and colleagues. As this involvement grows, we require higher levels of dependability. We expect computing systems to be trustable, reliable and secure from malicious attacks, to comply with all safety requirements, and to protect our privacy. Ensuring these properties will require more powerful methodologies and tools to design and implement dependable system in a cost-effective way.

# GLOSSARY AND ABBREVIATIONS

ASIC	<i>Application-Specific Integrated Circuits</i> are integrated circuits designed for a particular purpose, as opposed for general use in many different situations.
CAGR	<i>Compound annual growth rate</i> is a business and investing specific term for the smoothed annualized gain of an investment over a given time period.
Cloud computing	<i>Cloud computing</i> is a paradigm whereby computing power is abstracted as a virtual service over a network. Executed tasks are transparently distributed.
CMOS	Complementary Metal-Oxide-Semiconductor is a common technology for constructing integrated circuits. CMOS technology is used in microprocessors, microcontrollers, static RAM, and other digital logic circuits.
Declarative programming	<i>Declarative programming</i> is a programming paradigm that expresses the logic of a computation without describing its control flow. Many languages applying this style attempt to minimize or eliminate side effects by describing what the program should accomplish, rather than describing how to go about accomplishing it (the how is left up to the language's implementation). This is in contrast with imperative programming, in which algorithms are implemented in terms of explicit steps
EUV	<i>Extreme ultraviolet</i> lithography (also known as EUV or EUVL) is a next-generation lithography technology using an extreme ultraviolet (EUV) wavelength, currently expected to be 13.5 nm.
FDSOI	<i>Fully Depleted Silicon On Insulator</i> (MOSFETs). For a FDSOI MOSFET the sandwiched p-type film between the gate oxide (GOX) and buried oxide (BOX) is very thin so that the depletion region covers the whole film. In FDSOI the front gate (GOX) supports less depletion charges than the bulk transistors so an increase in inversion charges occurs resulting in higher switching speeds. Other drawbacks in bulk MOSFETs, like threshold voltage roll off, higher sub-threshold slop body effect, etc. are reduced in FDSOI since the source and drain electric fields can't interfere due to the BOX (adapted from Wikipedia)
Filter Bubble	A <i>filter bubble</i> is a situation in which a website algorithm selectively guesses what information a user would like to see based on information about the user (such as location, past click behavior and search history) and, as a result, users become separated from information that disagrees with their viewpoints, effectively isolating them in their own cultural or ideological bubbles.
FinFET	The term <i>FinFET</i> was coined by University of California, Berkeley researchers (Profs. Chenming Hu, Tsu-Jae King-Liu and Jeffrey Bokor) to describe a nonplanar, double-gate transistor built on an SOI substrate.... The distinguishing characteristic of the FinFET is that the conducting channel is wrapped by a thin silicon "fin", which forms the body of the device. In the technical literature, FinFET is used somewhat generically to describe any fin-based, multigate transistor architecture regardless of number of gates (from Wikipedia).
GPU	A <i>Graphics Processing Unit</i> refers to the processing units on video cards. In recent years, these have evolved into massively parallel execution engines for floating point vector operations, reaching performance peaks of several gigaflops.
HiPEAC	The European Network of Excellence on <i>High Performance and Embedded Architecture and Compilation</i> coordinates research, facilitates collaboration and networking, and stimulates commercialization in the areas of computer hardware and software research.
ICT	<i>Information &amp; Communication Technology</i> is a generic term used to refer to all areas of technology related to computing and telecommunications.
Imperative programming	<i>Imperative programming</i> is a programming paradigm that describes computation in terms of statements that change a program state. In much the same way that imperative mood in natural languages expresses commands to take action, imperative programs define sequences of commands for the computer to perform. The term is used in opposition to declarative programming, which expresses what the program should accomplish without prescribing how to do it in terms of sequences of actions to be taken.

Internet of Things	The Internet of Things (IoT) is a computing concept that describes a future where everyday physical objects will be connected to the Internet and will be able to identify themselves to other devices.
ISA	An <i>Instruction Set Architecture</i> is the definition of the machine instructions that can be executed by a particular family of processors.
JIT	<i>Just-In-Time</i> compilation is the method of compiling code from source or an intermediate representation at the time when it will execute. This allows for improved portability by generating the correct binary at execution time, when the final target platform is known. JIT compilation has been heavily leveraged in Java, Microsoft's C#, and OpenCL.
MEMS	Microelectromechanical systems is the technology of very small devices. MEMS are made up of components between 1 to 100 micrometres in size, and MEMS devices generally range in size from 20 micrometres to a millimetre.
NRE	<i>Non-Recurring Engineering</i> costs refer to one-time costs incurred for the design of a new chip, computer program or other creation, as opposed to marginal costs that are incurred per produced unit.
Programming model	A <i>programming model</i> is a collection of technologies and semantic rules that enable expressing algorithms in an efficient way. Often, such programming models are geared towards a particular application domain, such as parallel programming, real-time systems, image processing, ...
RFID	Radio-Frequency Identification is the use of a wireless non-contact system that uses radio-frequency electromagnetic fields to transfer data from a tag attached to an object, for the purposes of automatic identification and tracking.
SME	Small and Medium-sized Enterprise, a company of up to 250 employees.
SoC	A <i>System on Chip</i> refers to integrating all components required for the operation of an entire system, such as processors, memory, and radio, on a single chip.
Soft Errors	A <i>soft error</i> is a temporary wrong result, often caused by cosmic rays or temperature effects, not by a permanent failure of the circuit (which is called a hard error). With increasing integration the chances on soft errors are said to increase too.
STDP	<i>Spike-Timing-Dependent Plasticity</i> is a biological process that adjusts the strength of connections between neurons in the brain. The process adjusts the connection strengths based on the relative timing of a particular neuron's output and input action potentials (or spikes).
STREP	Specific Targeted Research Project – a type of European collaborative research and technology development project.
VLSI	Very-large-scale integration is the process of creating integrated circuits by combining thousands of transistors into a single chip.

# REFERENCES

- [Aitken11] M. Aitken, K. Flautner, and J. Goodacre, "High-Performance Multiprocessor System on Chip: Trends in Chip Architecture for the Mass Market," in *Multiprocessor System-on-Chip, Hardware Design and Tool Integration*, New York, NY: Springer New York, 2011, pp. 223-239.
- [Alibart10] F. Alibart et al. "An Organic Nanoparticle Transistor Behaving as a Biological Spiking Synapse," *Adv Functional Materials* 20.2, 2010.
- [Annun12] M. Annunziata, and P. Evans (2012) "Industrial Internet: pushing the boundaries of minds and machines", report, GE. November 26<sup>th</sup>, 2012, available at [http://www.ge.com/docs/chapters/Industrial\\_Internet.pdf](http://www.ge.com/docs/chapters/Industrial_Internet.pdf) retrieved on November 2012.
- [Arwu2012] <http://www.arwu.org> retrieved on November 2012
- [Bernstein] K. Bernstein et al. "Device and architecture outlook for beyond CMOS switches," *Proceedings of the IEEE* 98.12 (2010): 2169-2184.
- [Borgh10] J. Borghetti, G. Snider, P. Kuekes, J. Yang, D. Stewart, and R. Williams, "'Memristive' switches enable 'stateful' logic operations via material implication," *Nature*, vol. 464, no. 7290, pp. 873-876, Apr. 2010.
- [Cavin] R. Cavin et al. "A long-term view of research targets in nanoelectronics," *Journal of Nanoparticle Research* 7.6 (2005): 573-586.
- [Samp11] A. Sampson, W. Dietl, E. Fortuna, D. Gnanapragasam, L. Ceze, and D. Grossman. "EnerJ: Approximate Data Types for Safe and General Low-Power Computation," *PLDI 2011*, pp. 164-174, 2011.
- [Chua71] L. Chua, "Memristor-The missing circuit element," *Circuit Theory, IEEE Transactions on*, vol. 18, no. 5, pp. 507-519, 1971.
- [E10] "The data deluge: businesses, governments and society are only starting to tap its vast potential" (special report on "Data, data everywhere"). *The Economist*. Feb 2010. <http://www.economist.com/node/15579717> and <http://www.economist.com/node/1557443> retrieved on November 2012.
- [Esmat11] H. Esmailzadeh, E. Blem, R. St. Amant, K. Sankaralingam, and D. Burger, "Dark Silicon and the End of Multicore Scaling," *ISCA 2011*, pp. 365-376, 2011.
- [ESA] [http://www.esa.int/Our\\_Activities/Technology/Building\\_a\\_lunar\\_base\\_with\\_3D\\_printing](http://www.esa.int/Our_Activities/Technology/Building_a_lunar_base_with_3D_printing) retrieved on February 2013.
- [Hey03] A. Hey and A. Trefethen. "The Data Deluge: An e-Science Perspective," in F. Berman, G. Fox and A. Hey, Eds. *Grid Computing - Making the Global Infrastructure a Reality*, pp. 809-824. Wiley, 2003
- [Hey09] T. Hey, S. Tansley, K. Tolle (eds.) "The Fourth Paradigm: data-intensive scientific discovery," Microsoft Research book. 2009. <http://research.microsoft.com/en-us/collaboration/fourthparadigm/>
- [Imreo6] A. Imre, G. Csaba, L. Ji, A. Orlov, G. H. Bernstein, W. Porod, "Majority Logic Gate for Magnetic Quantum-Dot Cellular Automata", *Science*, Vol. 311 no. 5758 pp. 205-208, 13 January 2006.
- [KET] [http://ec.europa.eu/enterprise/sectors/ict/key\\_technologies/index\\_en.htm](http://ec.europa.eu/enterprise/sectors/ict/key_technologies/index_en.htm) retrieved on September 1<sup>st</sup>, 2011.
- [Kuzum11] D. Kuzum, R. Jeyasingh, B. Lee, P. Wong, "Nanoelectronic Programmable Synapses Based on Phase Change Materials for Brain-Inspired Computing," in *Nano Letters*, June 2011.
- [Lee12] H. D. Lee et al. "Integration of 4F2 selector-less crossbar array 2Mb ReRAM based on transition metal oxides for high density memory applications," in 2012 Symposium on VLSI Technology (VLSIT), 2012, pp. 151-152.
- [Manchester] <http://www.graphene.manchester.ac.uk/future/> retrieved on December 1<sup>st</sup>, 2012
- [Pacemaker] <http://techcrunch.com/2012/10/17/hacked-pacemakers-could-send-deadly-shocks/> retrieved on November 20<sup>th</sup>, 2012
- [Cheemo5] S. Cheemalavagu, P. Korkmaz, K. Palem, "Ultra low-energy computing via probabilistic algorithms and devices: CMOS device primitives and the energy-probability relationship," *International Conference on Solid State Devices*. Tokyo, pp. 2-4, 2004.
- [Partha12] <http://h30507www3.hp.com/t5/Innovation-HP-Labs/Electrons-for-compute-photons-for-communication-ions-for-storage/ba-p/115067> retrieved on November 2012.
- [Pershin11] Y. Pershin and M. Di Ventra, "Memory effects in complex materials and nanoscale systems," *Advances in Physics*, vol. 60, no. 2, p. 145, 2011.
- [Sampa] <http://sampa.cs.washington.edu/sampa/EnerJ> retrieved on November 25<sup>th</sup>, 2012
- [Schirber] <http://physics.aps.org/articles/v5/24> retrieved on December 1<sup>st</sup>, 2012
- [Snidero8] G. Snider, "Spike-timing-dependent learning in memristive nanodevices," in *Nanoscale Architectures*, 2008. NANOARCH 2008. *IEEE International Symposium on*, 2008, pp. 85-92.
- [Strukovo8] D. Strukov, G. Snider, D. Stewart, and R. Williams, "The missing memristor found," *Nature*, vol. 453, no. 7191, pp. 80-83, May 2008.
- [SyNAPSE] [https://www.ibm.com/smarterplanet/us/en/business\\_analytics/article/cognitive\\_computing.html](https://www.ibm.com/smarterplanet/us/en/business_analytics/article/cognitive_computing.html) retrieved on 1/9/2011.
- [TimeTech] [echland.time.com/2013/02/05/we-can-almost-print-new-organs-using-3d-stem-cells/](http://echland.time.com/2013/02/05/we-can-almost-print-new-organs-using-3d-stem-cells/) retrieved on February 2013.
- [Xiao9] Q. Xia et al., "Memristor-CMOS Hybrid Integrated Circuits for Reconfigurable Logic," *Nano Letters*, vol. 9, no. 10, pp. 3640-3645, Sep. 2009.



# SUMMARY

Energy has become the primary limiting factor in the development of all systems. This has led to parallel and heterogeneous devices, and the appearance of “dark silicon”. Developing energy-aware devices and automating the optimization of applications for power efficiency has become a necessity across all computing systems.

The traditional computing systems market is evolving towards *global-scale applications* that gather data from embedded systems and users, process it in large data centers, and provide customized, timely information to millions of users through their mobile devices. Tools have not yet caught up with this evolution towards global-scale applications.

Europe provides a strong embedded and low-power processor ecosystem with many companies in hardware and software development. However, Europe also suffers from a high degree of horizontal specialization, which makes it difficult for companies to amortize the costs of development across the product chain. The most successful global-scale applications are backed by global vertically-integrated companies offering a complete ecosystem.

## STRATEGIC AREAS FOR HORIZON 2020

**Embedded systems.** We need to rethink system architectures and programming models to optimize for energy, time constraints and safety, and develop techniques to support portability of critical and mixed-critical systems.

**Mobile systems.** More natural and immersive experiences will improve the utility and productivity of mobile devices. Further, we will need stronger guarantees of privacy and security.

**Data center computing.** We need to find ways to process “big data” without increasing cost or energy. We must integrate these developments with techniques to ensure security, privacy, and compliance, while providing large-scale reliability, availability, and serviceability.

## CROSS-CUTTING CHALLENGES FOR HORIZON 2020

**Energy efficiency.** We must address the challenges of programming parallel heterogeneous processors and optimizing data movement. We must also take advantage of the energy saving potential of new technologies.

**System complexity.** We need to develop tools and techniques to optimize for performance, ensure correct operation, all while operating “at-scale”. We need to investigate new integration techniques that enable high levels of integration and differentiation without the cost of cutting-edge fabrication.

**Dependability.** We need more powerful methodologies and tools to design and implement dependable systems in a cost-effective way.



[roadmap@hipeac.net](mailto:roadmap@hipeac.net)  
<http://www.HiPEAC.net/roadmap>